# An Optimized Hybrid Approach for Airline Ticket Pricing

Anala Lakshmi Anuroopa<sup>1</sup>, Maganti Venkatesh<sup>2</sup> <sup>1</sup> Aditya University, Surampalem, India, <sup>2</sup> Associate Professor, Aditya University, Surampalem, India,

Abstract. Air ticket pricing is volatile, as various dynamic factors driving change alter demand levels, seasonal variations, booking time, and competitive pricing. It indeed becomes confusing and difficult for the airlines as well as passengers because the prices of tickets often change erratically. With such an understanding, the potential to accurately estimate ticket prices has grown to become a vital revenue maximization tool in airlines while similarly helping consumers make educated purchase decisions. Therefore, considering that such price variations may be quite challenging and unpredictable, the demand is gigantic for robust prediction models of prices that will give reliable estimations under a myriad of market conditions. Five egression models to forecast air ticket prices are critically analyzed in this research work: Linear Regression, Polynomial Regression, Decision Tree Regression, Gradient Boosting Regression, and Gradient Boosting Regression optimized with Particle Swarm Optimization (PSO). The models selected have already established presence in predictive analytics and differ in complexity, interpretability, and computational efficiency. Performance of each model is measured and compared according to many metrics: Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Bias Deviation (MBD), R<sup>2</sup> Score, and Accuracy. Our results show that Gradient Boosting Regression optimized with Particle Swarm Optimization (PSO) achieves the highest predictive accuracy, and it is by far the most reliable choice for a ticket price forecast.

**Keywords:** Machine Learning, Airline Ticket Pricing, Predictive Modeling, Regression Models, Gradient Boosting, Decision Tree, Price Forecasting, Computational Efficiency, Overfitting, Particle Swarm Optimization (PSO),Real-Time Prediction

#### 1. Introduction

The airline business environment is rather dynamic and competitive where ticket prices change frequently and sometimes unpredictably during a particular day. Prices of airline tickets generally fluctuate according to a combination of variables such as the level of consumers, seasonalities, booking date, competitor's price strategies, and market conditions. This creates significant challenges both for the airlines and the consumers, since both must navigate an environment in which prices might change at any instance, preventing airlines from possibly unlocking greater returns and consumers from efficiently making a decision as to when to buy. Given the volatile nature of air ticket pricing, it is very important for airlines to implement methods which would enable their organizations to predict more accurate and consistent price of a ticket so that their decisions are better strategic.

Among various methods, use of predictive modeling using some techniques in machine learning (ML) has become one of the most viable solutions for predicting future ticket price. With historical data, these models can be used to spot trends, patterns, and correlations in pricing, thus providing an opportunity for airlines to predict ticket prices as a function of such factors as demand, booking windows, and competitor pricing. Such an ability could give airlines a huge competitive advantage, allowing them to optimize their pricing strategies, enhance revenue management, and inform consumers of more precise price predictions. With predictive modeling, airlines could better command real-time pricing management, an utmost need in making adjustments when market conditions change or competitors act to keep up.

This paper aims to provide further insight into the performances of common regression models applied in machine learning for airline ticket price predictions: Linear Regression, Polynomial Regression, Decision Tree Regression, Gradient Boosting Regression, and Gradient Boosting Regression optimized with Particle Swarm Optimization (PSO). The models vary in terms of complexity and interpretability, and indeed strengths and weaknesses in accuracy, computational efficiency, and ability to forecast. For example, a linear Regression model is one of the more interpretable models because it gives insights into linear relationships among the variables, while Gradient Boosting Regression optimized with Particle Swarm Optimization (PSO) provides higher predictive accuracy compared to some other models, with greater computational cost and a longer training time. The decision tree regression model requires a transparent approach that can easily be understood, but the performance depends on the depth and complexity of the decision tree in such models. On the other hand, Polynomial Regression can model nonlinear relationships but must be tuned with caution to avoid overfitting.

This study contributes valuable insights into the application of machine learning models in practice for predicting airline ticket prices. The comparison and evaluation of multiple regression techniques contribute to the continuing development of better, more efficient airline pricing strategies, enabling airlines to optimize revenues and serve consumers better within a rapidly changing market environment. The findings open avenues for future research that can develop accuracy and efficiency in price prediction models further enhanced perhaps with factors such as competitor pricing, customer behavior, and external economic indicators.

## **2. LITERATURE SURVEY**

Korkmaz [1] has applied many machine learning modeld on airline flight datasets and observed that GPR (Gaussian Process Regression) has achieved the high accuracy.

Deng [2] has analyzed various factors which will impact flight fares but failed to capturecomplex interactions between price determinants.

Zhang [3] has developed forecasting model on airline flight based on the passenger's historical data. That models can help optimize price strategies by past trends, predicts future patterns. This study has stated that time series methods like ARIMA and LSTM, it outperforms traditional regression models for airline price prediction

Ali et al.[4] has analyzed many regression models, which includes Linear Regression, Random Forest, Gradient Boosting and to forecast airline price dynamically. These observations states that ensemble methods will provide higher accuracy during the price fluctuations analysis, but traditional models underperformed on volatile prices.

Wang et al. [5] did the study on big data of airline pricing models with Apache Spark and machine learning models like Gradient Boosted Trees, Decision Trees etc, found that Random Forest models performs best as compared to traditional regression methods, when dealing with large-scale datasets.

Kumar et al. [6] was introduced a best price sensitive calculation framework, by using machine learning models like Poisson semi-parametric regression and the study states that improved forecasting of airline demand and pricing trends. The results indicated that hybrid Machine learning models can adapt more effectively on dynamic market conditions.

Yildirim et al. [7] has implemented deep learning models, Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks, for predicting airline prices. This experiment concludes that deep learning models can capture non-linear patterns in prices effectively than the traditional machine learning algorithms.

Sharma & Gupta [8] states that combination of AI in revenue management systems in airline price predictions, the framework has included the reinforcement learning based methodologies for adjusting the airline prices based on the past trends and patterns that demand competitor pricing.

Zhang et al. [9] studies the implementation Machine Learning models to predict airline forecating on prices. This focuses on the implementation of gradient boosting methods, specifically eXtreme Gradient Boosting (XGBoost), Light Gradient Boosting Machine (LightGBM), and CatBoost, optimized using Particle Swarm Optimization (PSO).

#### **Identified Gaps:**

- Most of the studies doesn't have real-time prediction and optimization, which will help in dynamic pricing systems.
- Critical external factors contribute more to enhance accuracy, but that were ignored.
- There were some computational challenges involved.
- Existing studies failed to propose solutions for generalization.

• Overfitting is the biggest concern particularly with complex models like Polynomial regression.

#### 2. 1 Need for the present work:

- Integrating interpretability to make the complex models more transparent.
- Evaluating linear and nonlinear approaches to improvise the accuracy with efficiency.
- Expanding the feature set to include market-driven variables, for stakeholders trust and adoption in industry practices.
- Analyzing computational bottlenecks in rea-time scenarios.

# **3. METHODOLOGY**

#### 3.1 Source of Data:

The dataset for this study consists of information regarding the prices of airline tickets. It introduces different influencing factors to the price of tickets. These influencing factors include attributes such as booking date, departure date, airlines, and other relevant variables that could influence ticket prices. A historical view of airline prices helps in building predictive models that can predict future ticket prices with regard to influencing factors.

#### 3.2 Data Preprocessing:

The preprocessing of the data is the process in which it is transformed to be suitable for model training. In this case, data processed through several stages of cleaning and transformation. Missing values in the dataset are replaced by the mean of its corresponding categorical variable. For example, categorical variables like airline names are encoded with the help of Label Encoding wherein each unique category is assigned a numerical value. This way, the categorical data will be interpreted by the machine learning algorithms appropriately and still provides integrity to the dataset. Other pre-processing involved feature scaling and normalization where applicable to ensure that the range of the data will not adversely affect the performance of the models.

#### 3.3 Architecture Diagram



# 3.4 Model Algorithms

## **Linear Regression**

Linear Regression is one of the most elementary and commonly applied statistical methods for modeling the relationships between variables. In the context of airline ticket price prediction, the goal is to model the price as a function of various factors, such as booking date, departure date, and airline. Linear Regression aims to find the line (or hyperplane in higher dimensions) that best fits the observed data.

Overview:

Cost Function (Mean Squared Error) In other words, the goal of Linear Regression is to minimize the difference between the predicted and the values that have been actually observed. This is calculated by the Mean Squared Error (MSE), which means the average squared difference in between the predictions and the true values. Optimization (Gradient Descent/Normal Equation): Once the cost function is defined, the model is trained by using methods like Gradient Descent or the Normal Equation in order to determine the best coefficients or weights assigned to the features in terms of minimizing the cost function.

$$y = \beta_0 + \beta_1 x + \epsilon \tag{1}$$

#### **Polynomial Regression**

Polynomial Regression is an extension of Linear Regression, using polynomial terms (higher order powers of the features) to model nonlinear interactions between the independent variables and the dependent variable. This is particularly useful for the cases in which the input variables' relationship with the output isn't strictly linear, but rather more curvaceous or of a higher degree.

Overview:

Creates polynomial features by transforming input features into a higher degree. A linear regression model is applied to these newly generated features to predict the target variable.

This method is very useful when there are relationships that are more complex and nonlinear than what a simple linear regression model can grasp.

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_n x^n + \epsilon$$
(2)

#### **Decision Tree Regression**

Decision Tree Regression is a non-linear approach for building a machine learning algorithm that constructs a tree-like model in order to predict an output variable. The algorithm recursively divides the feature space into smaller regions based on features, with each being associated with a prediction. For each decision tree node, the best feature and threshold are selected to the one minimizing a specific error criterion, for example, Mean Squared Error (MSE).

#### Overview:

Initialization: The algorithm starts with the whole dataset and works step-by-step to pick the best feature and threshold that splits the data into two subsets. Recursive Splitting is the process of selecting the best feature and threshold continues recursively, dividing data into smaller and smaller subsets at each node of the tree. It stops with some conditions: a minimum number of samples in a node or when further splits no longer improve the model's performance. Once the tree has been fully built, every node at the bottom represents the predicted value, in general the mean target value of the samples in that node. Decision Tree Regression is particularly useful because it can model complex, non-linear relationships and is highly interpretable, since the decision-making process can be visualized via the tree structure.

Variance Reduction = Var(T) - 
$$\left(\frac{|T_L|}{|T|} \cdot Var(T_L) + \frac{|T_R|}{|T|} \cdot Var(T_R)\right)$$
 (3)

#### **Gradient Boosting Regression**

Gradient Boosting Regression is a versatile ensemble learning method in which multiple decision trees are put together in a sequential manner; each tree is learned so as to correct the errors made by the previous ones. Unlike traditional decision tree models, Gradient Boosting improves the model iteratively, hence this is an advancement of the basic decision tree, and often such an improvement leads to superior predictive accuracy.

#### Overview:

The Gradient Boosting begins with an initial model which, most of the times, is a simple model that predicts the average of the target variable across all samples. In each iteration, a new decision tree is appended to the ensemble. It is trained in order to predict the residuals, or differences between the predicted values and the actual target values, of the previous model. After training a new tree, the predictions of all the trees in the ensemble are combined. The prediction made by the new model is a weighted sum of the predictions from all the individual trees. The process goes on till a specified number of trees is added, or the improvement in performance is negligible.

Every new tree tries to rectify the mistakes of the past ensemble and thus emphasizes the residuals. The technique is very much in use nowadays due to its higher accuracy and robustness, especially for complex data. This technique has become one of the most popular algorithms for regression tasks in machine learning.

Each of these regression models proposes an essentially different way of dealing with the intricacies of airline ticket price prediction. While the approach in Linear Regression is simple and straightforward to interpret, more complex approaches such as Polynomial Regression and Decision Tree Regression are available to capture nonlinear relations. Gradient Boosting Regression, although computationally more expensive, often offers far better predicting performance, making it a fantastic approach for real-world applications such as price forecasts.

$$F_t(x) = F_{t-1}(x) + \nu \cdot h_t(x)$$
 (4)

# Gradient Boosting Regression optimized with Particle Swarm Optimization (PSO):

The search space for hyperparameters (e.g., learning rate, tree depth, number of trees) is defined. Each particle in the swarm is initialized with random values within this space.

1. The distance between its current position and its personal best (pbest), which encourages exploration of its own best-performing solutions.

2. The distance between its current position and the global best (gbest), which encourages convergence towards the overall best solution found by the swarm.

$$v_i(t+1) = w \cdot v_i(t) + c_1 \cdot r_1 \cdot (pbest_i - x_i(t)) + c_2 \cdot r_2 \cdot (gbest - x_i(t))$$
(5)

# **4. RESULTS AND DISCUSSION**

#### 4.1 Model Performance Comparison

Model	MAE
Linear Regression	4020.01
Polynomial Regression	6044.15
Decision Tree Regression	4335.32
Gradient Boosting Regression	3063.31
Gradient Boosting Regression (PSO)	3188.04

#### Table 1: Performance metrics for MAE



# Fig 2: MAE vs Model

Model	MSE	
Linear Regression	30818586.3	
Polynomial Regression	64533992.18	
Decision Tree Regression	48184243.37	
Gradient Boosting Regression	22507867.39	
Gradient Boosting Regression (PSO)	26736342.52	





Fig 3: MSE vs Model



#### Table 3: Performance metrics for RMSE



 Table 4: Performance metrics for MBD

Model		MBD
Linear Regre	ession	258.6
Polynomial Reg	gression	46.9
Decision Tree R	egression	565.39
Gradient Boosting	Regression	265.32
Gradient Boosting Reg	gression (PSO)	433.4
MBD	vs. Model	
	565.39	433.4
259.6	265.3	2



MBD

Model Fig 5: MBD vs Model



 Table 5: Performance metrics for R<sup>2</sup> Score



Table 6: Performance metrics for Accuracy

Model	Accuracy
Linear Regression	0.8
Polynomial Regression	0.83
Decision Tree Regression	0.87
Gradient Boosting Regression	0.9
Gradient Boosting Regression (PSO)	0.94



Fig 7: Accuracy vs Model

**PAGE NO : 101** 

# 5. DISCUSSION

This study took a deep dive into how well different regression models can predict airline ticket prices by looking at various factors like flight date, airline, class, source, destination, timings, duration, stops, and more. The models we examined included Linear Regression, Polynomial Regression, Decision Tree Regression, Gradient Boosting Regression, and a special version of Gradient Boosting that uses Particle Swarm Optimization (PSO). To measure how well these models performed, we used several key metrics: Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Bias Deviation (MBD), R<sup>2</sup> Score, and Accuracy.

Among the models, Gradient Boosting Regression stood out with the lowest MAE (3063.31) and RMSE (4744.25), along with the highest R<sup>2</sup> Score (0.95). This indicates that it did an excellent job of fitting the data and keeping prediction errors to a minimum. When we added PSO optimization to Gradient Boosting Regression, it achieved the highest Accuracy (94%). This shows that the optimization really boosted our confidence in the classifications, even though it was slightly behind Gradient Boosting in terms of RMSE and MAE. Still, it maintained a strong R<sup>2</sup> Score of 0.94, proving its solid predictive capabilities.

Linear Regression provided a decent starting point with an MAE of 4020.01 and an  $R^2$  Score of 0.93, making it a good choice for simpler situations, though it didn't quite measure up to the more advanced tree-based models. On the other hand, Polynomial Regression struggled the most, showing the highest MSE (64533992.18) and the lowest  $R^2$  Score (0.86). Its complexity likely led to overfitting, which hurt its ability to generalize well.

Decision Tree Regression had a balanced performance, with a moderate MAE (4335.32) and Accuracy (87%), but it still fell short compared to the Gradient Boosting models in all the important areas.

In summary, while all the models managed to learn from the dataset and deliver reasonable predictions, Gradient Boosting Regression clearly emerged as the most dependable option, striking the best balance between low error rates and high  $R^2$  scores.

# 6. CONCLUSION

The experimental results show that the PSO-enhanced Gradient Boosting Regressor outperformed all major evaluation metrics. With an impressive accuracy of 94%, an RMSE of 5170.72, and a solid R<sup>2</sup> score of 0.94, this model proved to be both precise and capable of generalizing well. By incorporating Particle Swarm Optimization (PSO)—a technique inspired by the social behaviors of birds and fish—this model achieved effective hyperparameter tuning that traditional Gradient Descent methods couldn't

match. This hybrid approach helped avoid local minima traps and sped up convergence rates, leading to a more stable and optimized model.

Additionally, the comparative analysis shed light on the limitations and tradeoffs of each model. For example, while Polynomial Regression initially seemed to fit non-linear trends well, it ended up overfitting and showing high variance with unseen data. Decision Tree Regression, while easy to interpret and quick, had moderate prediction errors because it tended to segment data too finely. In contrast, Gradient Boosting and its PSO variant excelled due to their ensemble nature and iterative refinement processes.

In a broader sense, this study adds to the growing field of applied machine learning for dynamic pricing, providing not just a high-accuracy solution but also a flexible framework for similar regression tasks. The proposed model is scalable, explainable, and modular, making it suitable for various industries like e-commerce, hospitality, and logistics, where accurate price prediction is crucial.

From an academic standpoint, this research connects the dots between machine learning, optimization algorithms, and specific problem-solving in different domains. It highlights the significance of selecting the right model architectures, tuning strategies, and evaluation methods that align with the dataset's characteristics and the business context. It also illustrates the value of interdisciplinary approaches in tackling complex challenges.

#### 7. FUTURE WORK

For improving regression models, including better performance on airline ticket price prediction, the following could be significant areas to focus on:

**7.1 Temporal and Seasonal Modeling:** Time-series forecasting techniques and recurrent neural network architectures like LSTM (Long Short-Term Memory) can be incorporated to model the temporal dependencies and seasonal variations in ticket prices more accurately, particularly for long-term predictions.

**7.2 Incorporation of Economic and Competitive Factors:** The current type of models depicts past records concerning prices of airline tickets. Future work can benefit by adding other external factors affecting prices in the broader set. These could be economic factors (inflation, fuel prices, growth in GDP) and competitive factors (competitor airline pricing strategies). This would add strength to the model, making it more responsive to changes in price when changes in the market occur. It would call for advanced feature engineering and data gathering from a number of different sources to compile a comprehensive dataset.

**7.3 Explainable AI Approaches Towards Interpretable Models:** As ensemble models like Gradient Boosting also start to become complex, it will be difficult to understand why the model is making a certain prediction. Exploring integrations of XAI approaches into such complex models would be an interesting direction to explore further. One example of XAI methodology is SHAP values or LIME (Local Interpretable Model-agnostic Explanations), which could shed light on how specific features contribute to model predictions. This would be particularly useful in industries where transparency in decision-making is crucial, allowing analysts to interpret and justify model outputs in practical applications.By focusing on these areas, future developments can make these models more efficient, adaptable, and interpretable in real-world applications, providing more actionable insights for airline pricing strategies.

# 8. REFERENCES

- Korkmaz, H. (2024). Prediction of Airline Ticket Price Using Machine Learning Method. Journal of Transportation and Logistics, 9(2), 205-218. https://doi.org/10.26650/JTL.2024.1486696
- [2] Deng, T. (2024). International flight fare prediction and analysis of factors impacting flight fare. Theoretical and Natural Science, 31(1), 329-335. https://doi.org/10.54254/2753-8818/31/20241079
- [3] Zhang, Y., Li, X., & Wang, H. (2024). Demand Forecasting Model for Airline Flights Based on Historical Passenger Flow Data. Applied Sciences, 14(23), 11413. https://www.mdpi.com/2076-3417/14/23/11413
- [4] Ali, M., Khan, R., & Sadiq, M. (2024). Predicting Optimal Airline Ticket Prices using Regression Models. Journal of Air Transport Research, Elsevier.

https://www.sciencedirect.com/science/article/pii/S1110016825001358

- [5] Wang, H., Zhao, T., & Liu, J. (2024). Using Spark Machine Learning Models to Perform Predictive Analysis on Flight Ticket Pricing Data. arXiv preprint arXiv:2310.07787. https://arxiv.org/abs/2310.07787
- [6] Kumar, A., Singh, P., & Roy, D. (2024). Machine Learning Based Framework for Robust Price-Sensitivity Estimation with Application to Airline Pricing. arXiv preprint arXiv:2205.01875. https://arxiv.org/abs/2205.01875
- [7] Yildirim, S., Akcay, H., & Celik, E. (2024). Deep Learning-Based Airline Pricing Prediction: A CNN-LSTM Approach. Journal of Transport and Logistics, 9(2), 45-62.
  - https://dergipark.org.tr/en/pub/jtl/issue/89876/1486696
- [8] Sharma, V., & Gupta, R. (2024). AI-Driven Revenue Management for Airline Pricing Optimization. Machine Intelligence Research, 12(1), 204-219.

http://machineintelligenceresearchs.com/index.php/mir/article/view/204

[9] Zhang, X., Wang, Y., & Liu, Z. (2023). Smart prediction of liquefactioninduced lateral spreading. Machine Intelligence Research, 21(4), 652– 669. https://doi.org/10.1007/s11633-023-1442-8