

Week 9 Milestone: Final Presentation Outline

Flight Price Predictor

Business case and problem statement:

Travelers often struggle with knowing when to book flights to secure the lowest prices. Airline pricing algorithms are complex, and last-minute purchases often result in higher costs. By using time series forecasting techniques and machine learning models, we aim to:

- Identify price trends over time
- Predict future ticket prices based on historical data
- Provide insights for consumers to make cost-effective purchasing decisions

Modeling hypothesis and assumptions:

- Model Hypotheses:
 - Flight prices generally increase as departure date approaches
 - Price fluctuations follow discernible patterns based on time to departure
 - Historical price data can predict future price trends
- Key Assumptions:
 - **Non-stationary data:** Means and variance change over time.
 - **Seasonal trends:** Holidays, weekends, and major events affect prices.
 - **Consistent relationship:** Time-to-departure influences price predictability.
 - **External factors:** Events and demand fluctuations impact pricing but can be modeled.
 - **Data availability:** Sufficient historical data exists for reliable predictions.

Data description and properties:

Dataset: 82 million flight itineraries with 46 features related to the flight itineraries.

Key Variables:

Variables	Description
totalFare (Target)	The final ticket price for a given flight departure.
search_date	The date the price is being checked.
departure_date	Date and time of departure.
seatsRemaining	Number of available seats at time of booking.
isRefundable	Indicates whether the ticket is refundable.
totalTravelDistance	Distance traveled (some missing values).

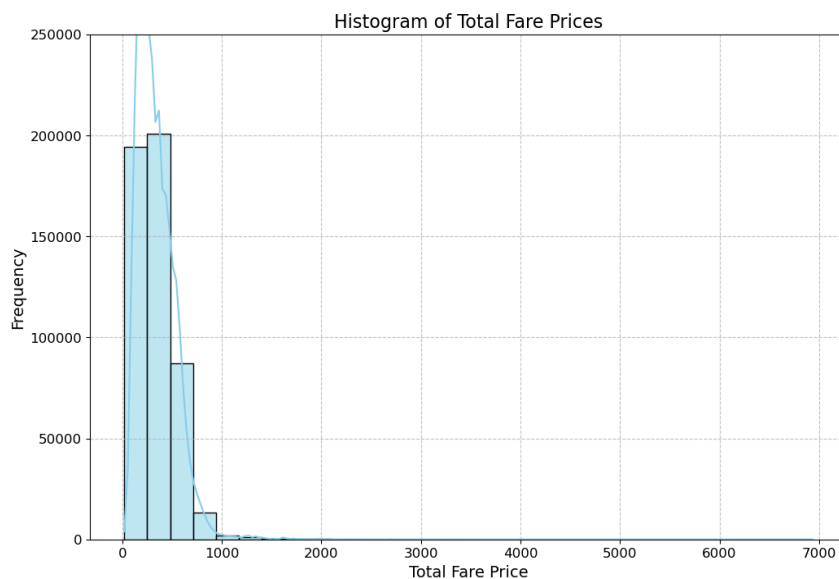
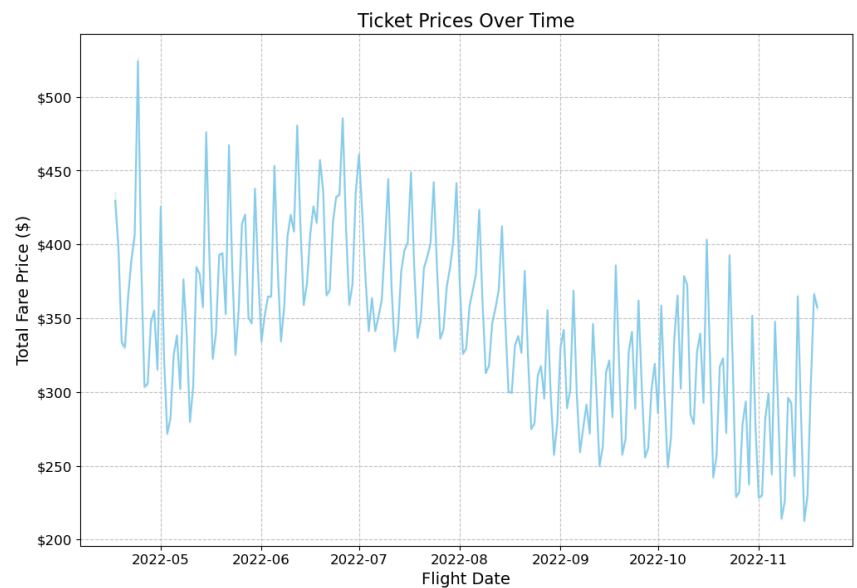
Data Challenges & Processing:

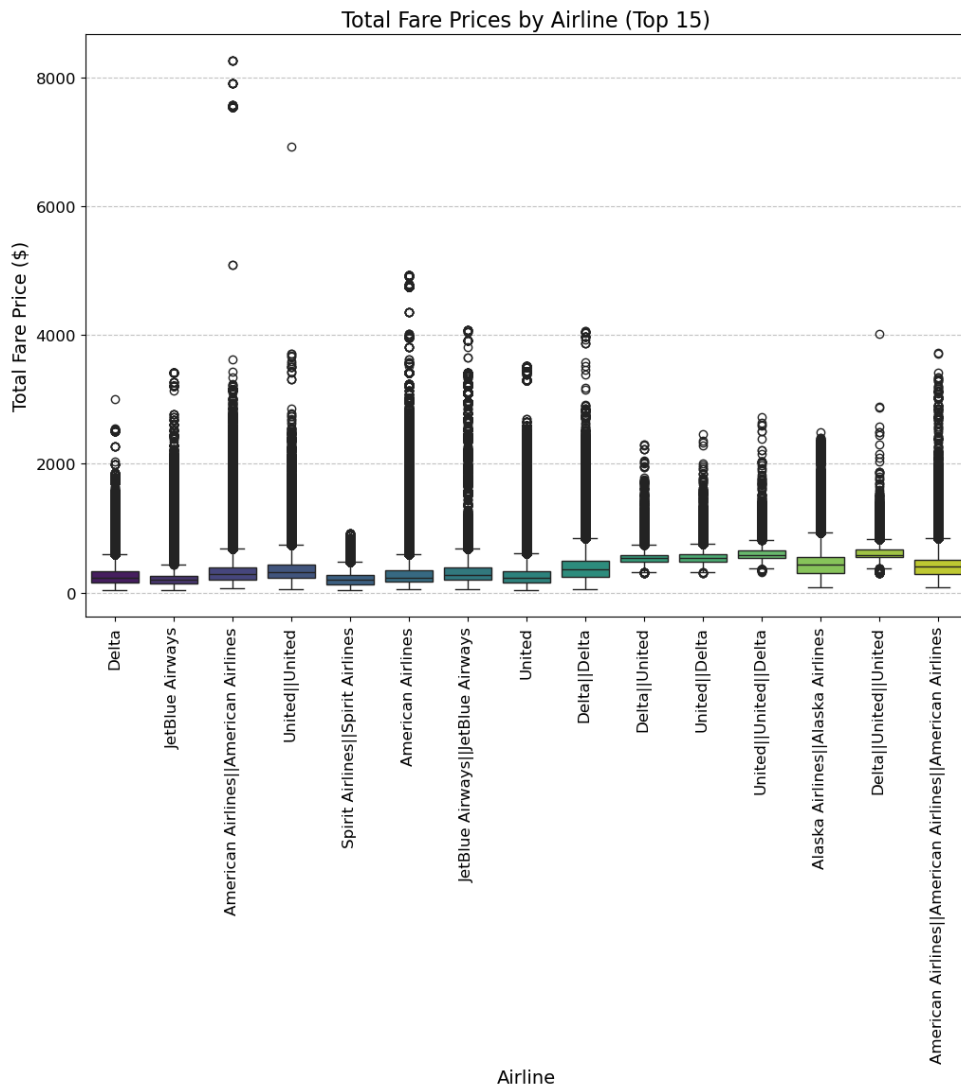
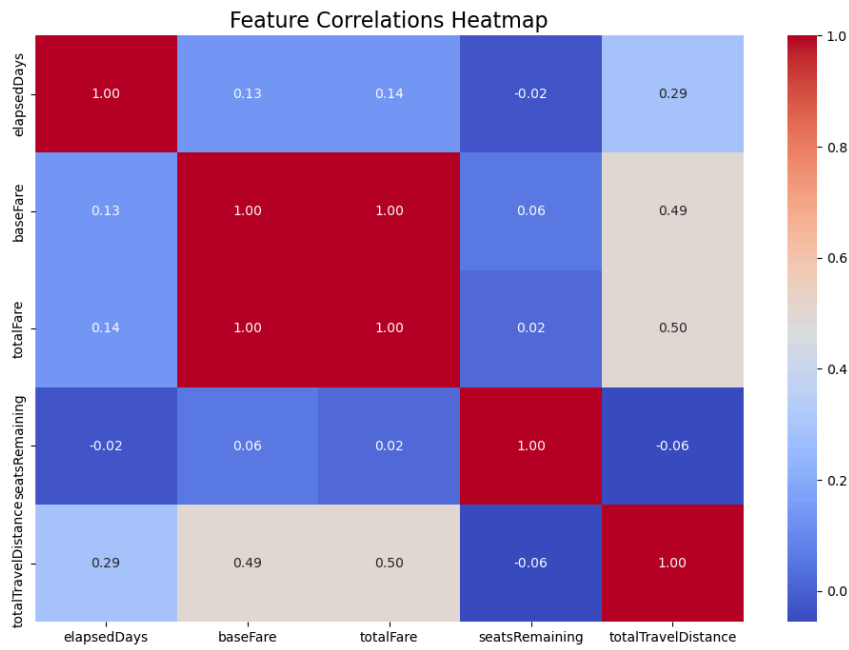
- Missing values in totalTravelDistance → Imputed using median-based techniques.
- Categorical variables (e.g., airport codes, airlines) → Encoded using one-hot encoding.
- Non-stationary trends in ticket prices → Feature transformations were applied.

Key Data Insights

- Flight prices increase closer to departure, but the rate varies based on seasonality.
- Some airlines exhibit more volatile pricing trends than others.
- The number of available seats influences price, lower availability leads to higher fares.

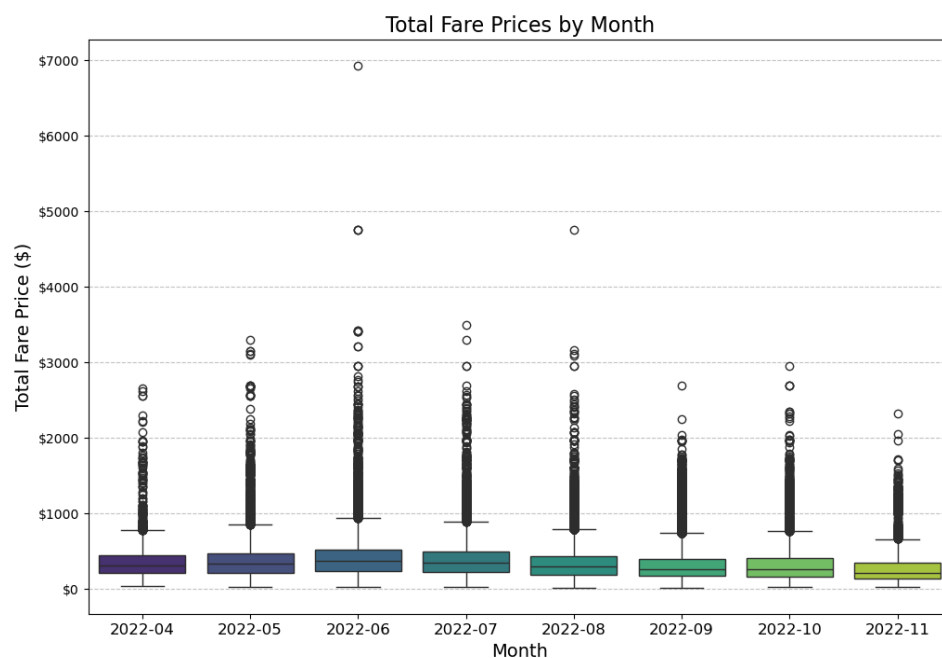
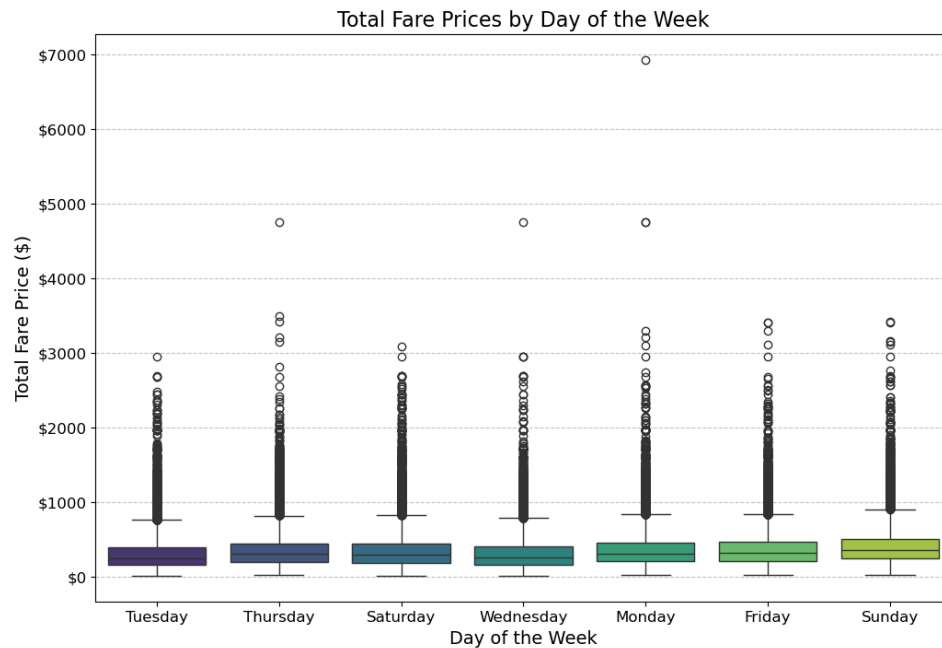
Data processing and Exploratory Data Analysis:

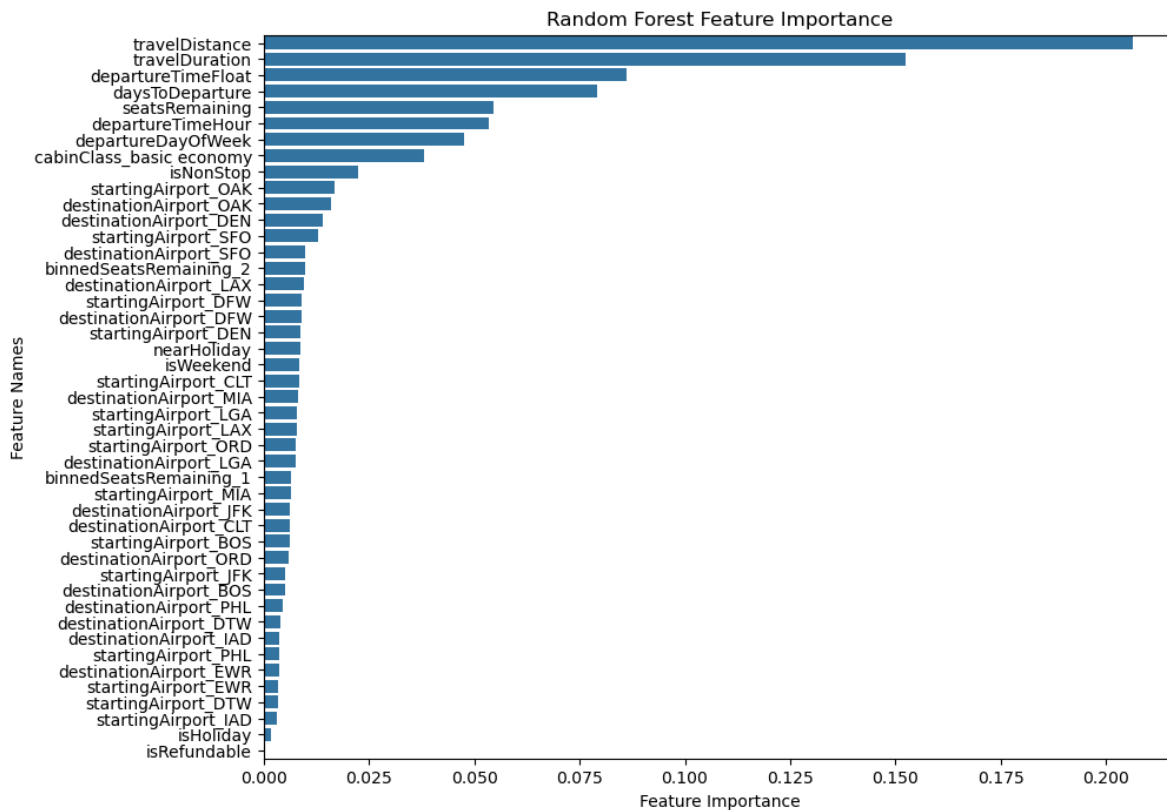




Feature Engineering:

1. *Derived Features:*
 - Binned seatsRemaining into categories (low, medium, high availability).
 - Extracted temporal features (e.g., month, weekday, hour of departure).
 - Created a holiday indicator → Flights near holidays flagged for price impact.
2. *Encoding Categorical Variables:*
 - One-hot encoding for airports and airlines.
3. *Missing Data Handling:*
 - Imputation of missing totalTravelDistance values.





Proposed modeling approaches:

Models Tested:

Traditional Time Series Models

- ARIMA
- SARIMA
- Prophet (did not generalize well)

Machine Learning Models

- Decision Tree (Baseline)
- Random Forest (Best-performing model)
- XGBoost (Close second in performance)

Deep Learning Model

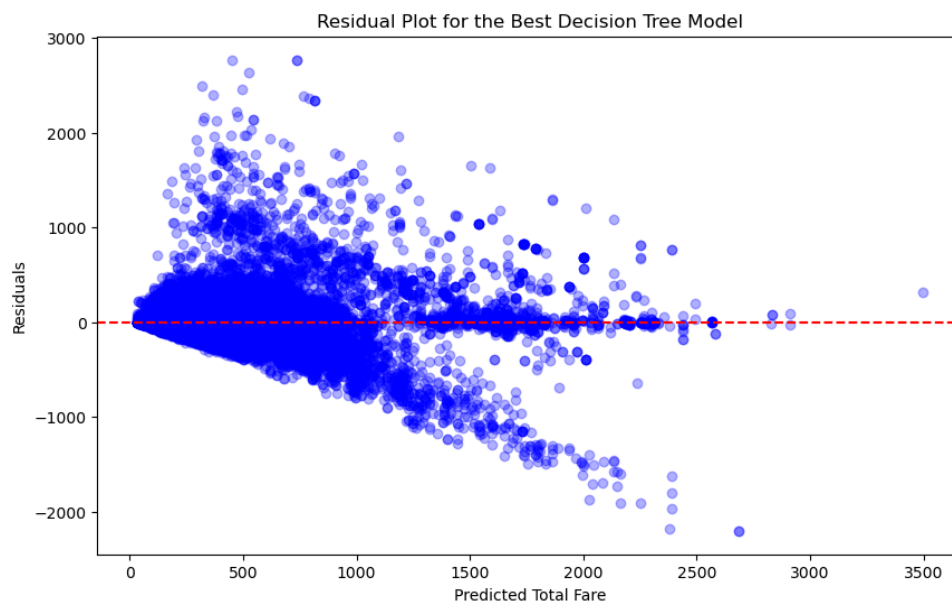
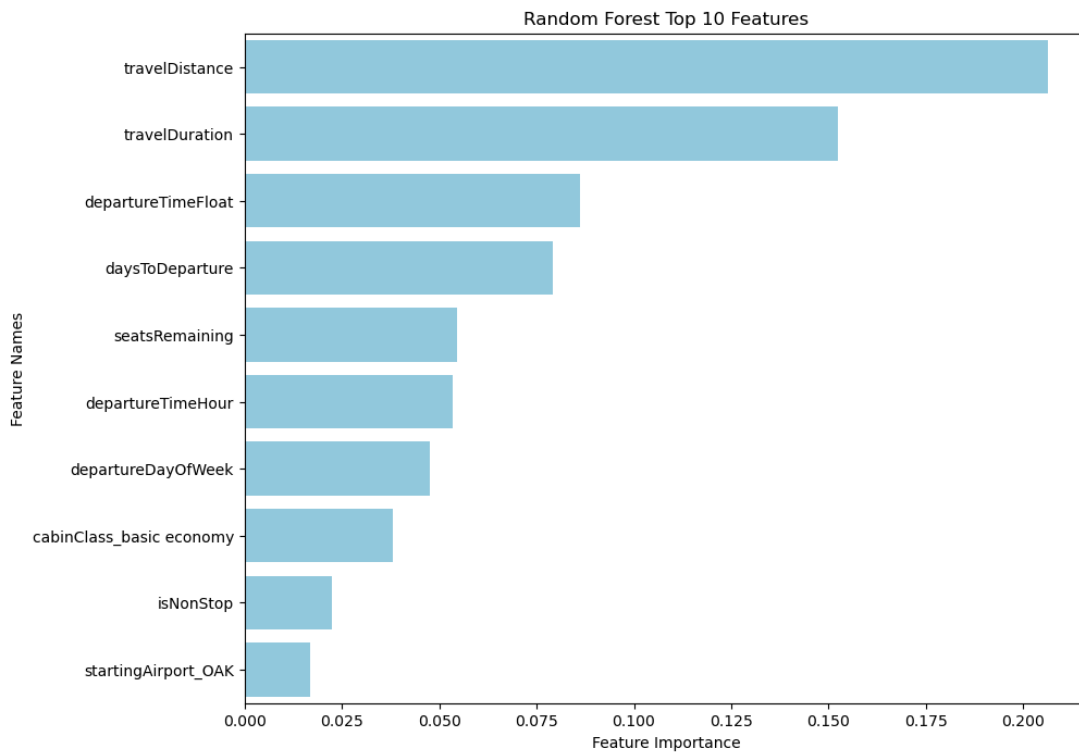
- RNN LSTM (Struggled due to price variance)

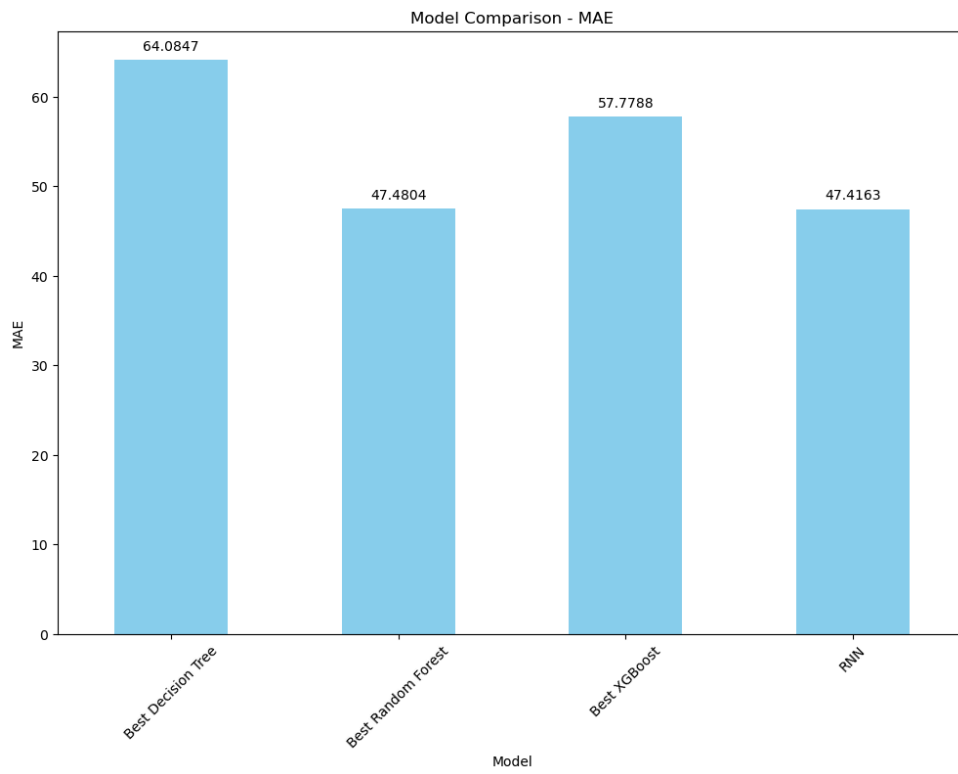
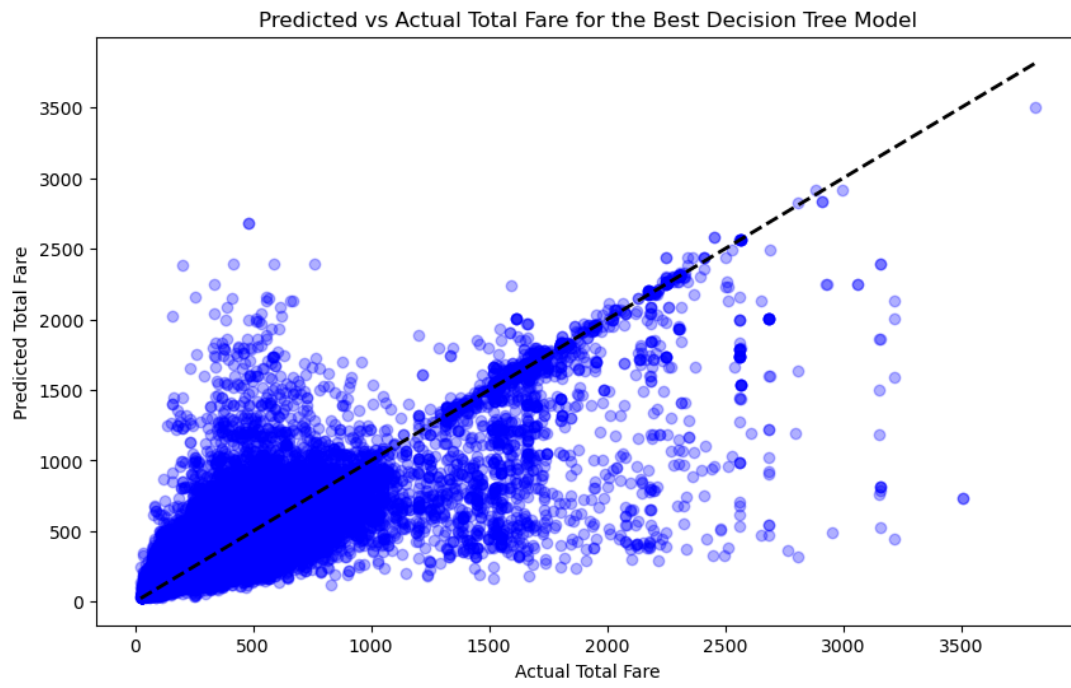
Selected model results with justifications and tradeoffs:

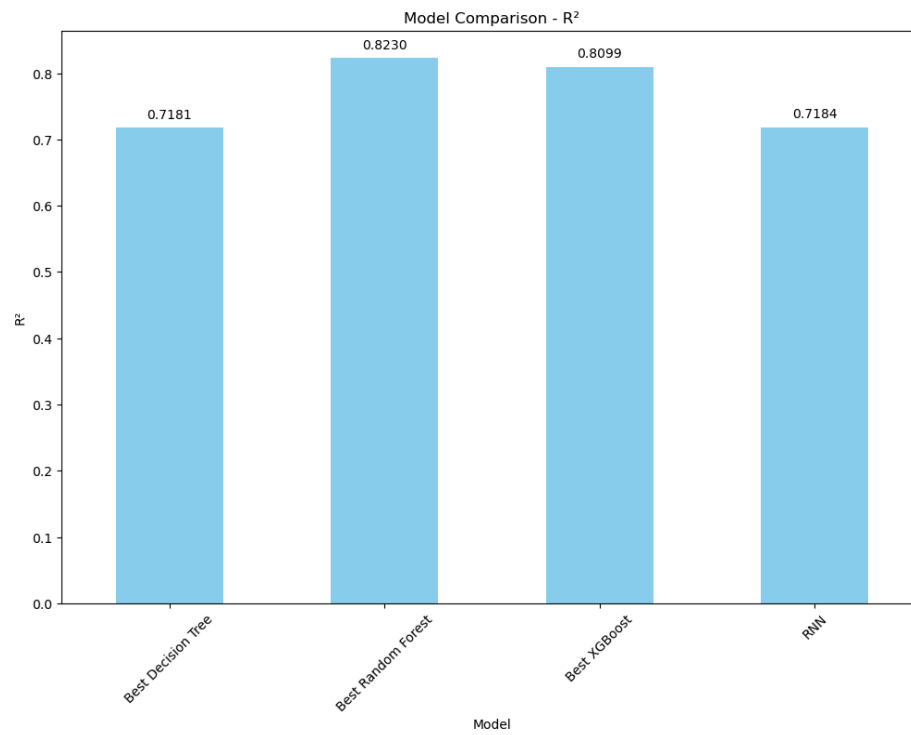
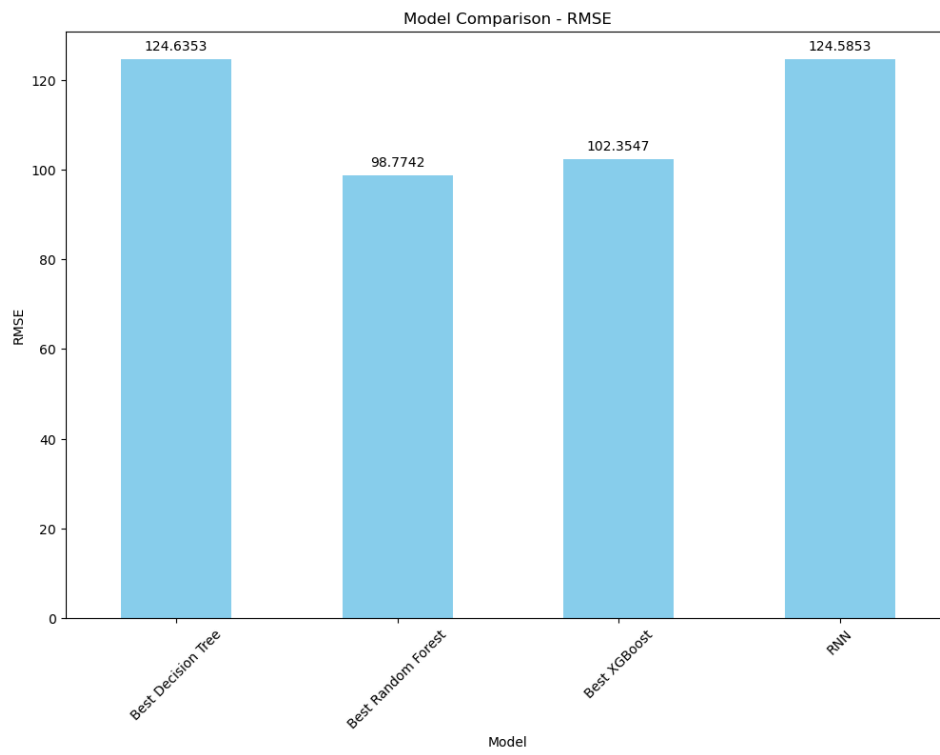
Model	BIC
ARIMA	55,701.33
SARIMA	52,994.49

Model	MAE	RMSE	R ² Score
Prophet	23,447.19	149.89	
XG Boost	56.60	99.70	0.82
Decision Tree	63.81	124.56	0.72
Random Forest	44.87	93.76	0.84
RNN LSTM	85.46	166.38	0.42

- Random Forest was the best-performing model, with an R² score of 0.84 and the lowest RMSE.
- XGBoost was a strong contender, showing slightly lower accuracy than Random Forest.
- Decision Trees had higher error rates due to overfitting on training data.
- RNN LSTM struggled with high variance, likely due to limited sequence dependencies in flight pricing data.







Insights/Recommendations & Future work

Insights from the Analysis:

- Tree-based models (Random Forest, XGBoost) outperform deep learning models for flight price prediction.
- Feature engineering significantly improved model accuracy, particularly holiday effects and seat availability binning.
- Flight prices are highly non-stationary, requiring advanced techniques for accurate forecasting.

Recommendations:

- Real-time prediction system: Deploy a model API that continuously updates with new flight data.
- External data sources: Incorporate weather conditions, real-time seat availability, and airline promotions to improve accuracy.
- Hyperparameter tuning: Further refine model parameters for even lower error rates.
- Transformer-based models: Test newer deep learning approaches (e.g., Time Series Transformers) to improve sequence-based learning.

Future Research Directions:

- Expanding the dataset to include international flights and budget airline fares.
- Exploring ensemble methods that combine Random Forest and XGBoost for improved predictions.
- Investigating reinforcement learning for dynamic pricing recommendations.