

# Real Time Hotel Booking Demand Optimization

Dr. D. Kavitha<sup>1</sup>, Aditya Kumar Singh<sup>2</sup>, Sakshi Chauhan<sup>3</sup>

Assistant Professor (Sr.), School of Computer Science and Engineering, Vellore Institute of Technology, Chennai<sup>1</sup>  
Students, School of Computer Science and Engineering, Vellore Institute of Technology, Chennai<sup>2,3</sup>  
kavitha.d@vit.ac.in, adityakumar.singh2020a@vitstudent.ac.in, sakshi.chauhan2020@vitstudent.ac.in

**Abstract:** *The significance of having proper revenue management and convenient operations in hotels rely on accurate daily demand. The main challenges that are being tackled conventionally are the forecasting of the number of cancellations and the calculation of the Average Daily Rate (ADR). One of the distinct characteristics of the hotel industry being so unstable is to adopt proactive strategies by dealing with a lot of external changes such as pandemics, disasters caused by nature, and economic fluctuations around the world. Booking cancellations are instrumental in helping hotel managers to optimise resources and inventory while ADR forecasting offers them equally essential projections concerning anticipated profit or loss margins. This research relies on several data modelling steps, including time series aggregation and decision merging, which are later followed by decomposition and model selection. SARIMAX and LSTM models are adopted for future traffic flow modelling, which demonstrates better forecasting performances. Binary classification is employed for feature engineering techniques together with model selection methods. Binary Classification is performed with a number of experiments with machine learning algorithms, AdaBoost turned out to be the best model which surpassed CART, KNN, Random Forest, Gradient boosting algorithms, and Light Gradient Boosting algorithms. The results of this are of great help to the hotel management for taking better decisions which are connected to the changing situations of the market*

**Keywords:** Daily Rate (ADR), Forecasting, Booking cancellations, SARIMAX, Long short-term memory (LSTM), Binary classification, K-nearest neighbours algorithm (KNN), Random Forest, AdaBoost, Gradient boosting algorithms, and Light Gradient Boosting algorithms

## I. INTRODUCTION

With dynamic and unpredictable business conditions, the hospitality industry needs to have an around-the-clock grasp of the actual demand in real time. This will help them actively manage their revenues and day-to-day activities within the hotel. This research study is proposed to be accomplished by two key strategies which are; predicting future cancellation of the reservations and the calculation of Average Daily Rate (ADR). Because of its complexity, not to mention the fact that its performance is greatly impacted by external factors like pandemics, natural disasters, and swings in the world economy, the hospitality business tends to be vulnerable to uncertainty [1]. The operational efficiency of a hotel is highly dependent on predictability of cancellations of reservation days; hence, the resource allocation can also be optimised for accuracy too. In addition, ADR forecasting is an analytical tool that helps management predict the probable profit margins and losses in hotel operation. Therefore, the research helps to understand the significance and efficiency of employing imaginative solutions to a wide range of problems, especially when the hospitality industry shows signs of declining. Seasonality is one of the main point factors considered to spot patterns and tendencies existing in data while forecasting hotel demands. Seasonality is a matter of these repeat starting and stopping cycles within set timeframes that could greatly affect the demand for hotel bookings [6]. Owners and managers of hotels need to be well-informed about seasonal events, holidays and festivals which definitely lead to high demand for these services and they must prepare in advance to meet these increased numbers of customers. This next juncture of time will finish by presenting the result of the periodical study, which is a more accurate Average Daily Rate (ADR) forecast and the proactive action concerning pricing schemes, resource allocation, service organisation, and planning [8]. Firstly, the decomposition method of the data is applied, and afterwards, diligent model selection is employed, that means systematic data verification and collation of time series through the process of merging and aggregation is done. SARIMAX and LSTM models for ADR forecasting in this study show a superior performance in contrast to the traditional methods. Moreover, the forecast of room cancellations is performed by applying binary

classification methods to enhance the feature engineering and model assessment [14]. In this research the comparison of different machine learning models in the thorough examination of which one works the best is done. The best model among all is AdaBoost which outperforms other models including CART, KNN, Random Forest, Gradient Boosting, and Light Gradient Boosting algorithms. These outcomes are the foundation for the decision support tool of the hotel management, where the latter can come up with appropriate measures to meet the shifts in the market.

## **II. LITERATURE REVIEW**

There are various factors that need to be considered by hospitality industry managers while budgeting and forecasting revenue. Through Rajopadhye et al., in 2021's proposition of unconstrained revenue management, forecasting of room demand found to suffice the incorporation of cancellation of reservations too according to the research. The researchers have to rely on true booking data also of uncompleted nights to make an informed and responsive forecast. Considering the combination of occupancy and arrivals on a weekly basis. M. Chattopadhyay and S. K. Mitra (2019) [17] used two year's daily arrival data, which was grouped into whole categories of the rate and length of stay to measure their forecasting accuracy at four different scales of aggregation. The disaggregation forecaster outperformed by a considerable degree all the other forecasters providing an aggregated forecast [2]. Thus, they suggested telling hotels to unify review data for guests together, categorising the length of stay and a room price range. Because of the low probability there will be any data about rate category and room type, hotels will need data by either rate category or room type. Karim and Weatherford (2023) point out, daily arrivals data brings a whole new level of precision to forecasting. Finally, they built a helpful forecast model employing Choice Hotels' performance data, and exploration and analysis utilising data obtained from Marriott Hotels helped confirm their assumptions [15]. Chen and Kachani (2017) suggested a forward model optimization of the network that is based on the links. This classical model involved the employment of a blend of traditional and modern techniques such as pick-up as well as simple exponential smoothing and the use of the combination of two for forecasting both the arrivals and occupancy data [11].

### **2.1 Real - estate economic trends**

Groundbreaking studies on the real estate sphere are being done through different approaches, each reviewing how the economy in this sector is fluctuating. According to the research of Quigley (2022), housing prices are lagged values and contain correlated components that in turn add to explanatory power, thus the housing cycle often shows robust price persistence. The survey data were established from the 41 MSA (metropolitan statistical areas) in Europe [3]. The research has led to the strong observation of a greater rate of global gross domestic product growth alongside the returns of commercial real estate globally implying that there is a direct association influenced by the global economy with real estate performance (Lieser & Groh et al., 2020) [4]. As an asset class, real estate is pretty close to the economic cycle. The economic health reflects the shifts in consumer behaviour, opportunities & threats in supply-demand systems and the investment preferences. For all these groups comprising investors, law makers, developers and practitioners, it is crucial to learn the financial structures of the real estate market.

### **2.2 Housing market dynamics**

In the context of the housing market evaluated by real estate literature, economics studies focus on it. Damary Glaeser and Gyourko's research indicates how restrictive housing supply, land-use laws, and population redistributions affect affordability and the price dynamics of housing. Moreover, it is emphasised by James M. Poterba's (2015) research, that demand of housing is highly susceptible to fluctuations of interest rates, as well as of family income, which would lead to decrease in housing expenditures [7].

### **2.3 Factors affecting demand for lodging**

Forecast of demand in lodging is impacted by the seasonality factor. This would include research of trend determinants and evaluation of some economical trends at different points in time (Parul Gupta, 2012) [8]. GDP mostly was seen as a very strong or the strongest predictor of hotel demand when it comes to looking at patterns in lodging demand in the European markets on the whole. The other areas in hospitality literature that take the view of demand forecasting are the studies on revenue management.

### 2.4 Machine learning methods in hotel forecasting

In order to improve the KDD program that would be employed in the ML approach of the hotel’s forecasting, the authors (Urraca et al., 2015) decided to use the genetic algorithms [5]. This test-set of the previous year's fund house performance from the period between January and July, was prepared. The hotel receives guests from different parts of Spain and also international tourists; as a result, reservation data was very essential. The collected data was from the hotel located in the northern province of Rioja. The data was utilised to train and validate the models for the years of 6. Moreover, the experts together picked up the best 119 indicators for attributes such as common traditions, local and national events, weather data, and socio-economic elements. After the discovery and visualisation of all the possible indicator dependencies by means of scatterplots and scatterplot matrices, these indicators were narrowed down to 22 ones. In forecasting accuracy models root mean squared error (RMSE) was used, GDF (generalised degrees of freedom) was applied to select the less complex models [12].

### III. DATA AND VARIABLES

The dataset is collected as a csv file with several columns. All hotels of this chosen data set are within Lisbon city, Portugal. The dataset primarily consists of two types of hotels: city hotels and resort hotels, where the perfect classification is considered for each class separately [9]. The database collection contains various reservation forms for both city and resort hotels, and the set of details covered by such information are quite [10]. The information to be integrated includes the specific time when reservations are secured, the duration of each reservation, and the number of adults, kids, and/or children associated with such reservation. The data comprises other variables such as parking availability and others in addition to them. Concurrently, the features present here produce information about the demand of real-time hotel bookings optimization, which is the primary target to attain.

### IV. RESEARCH METHODOLOGY

#### 4.1 Data collection

In the phase of data collection, first the time series demand data were collected directly from the hotel management software. Subsequently, the Hotel booking demand dataset was acquired, consolidated within a single CSV file featuring multiple columns. The dataset primarily comprises two distinct types of hotels: urban hotels and resorts, each of them carefully marked-off into the corresponding category.

#### 4.2 Data pre-processing

The first phase objective is to standardise and refine the hotel demand dataset. The dataset is processed with the use of descriptive analytics tools which help in finding the patterns or relationships in the data. Following this, feature engineering methods are applied in congruence with the convenient methods that handle missing data [13].

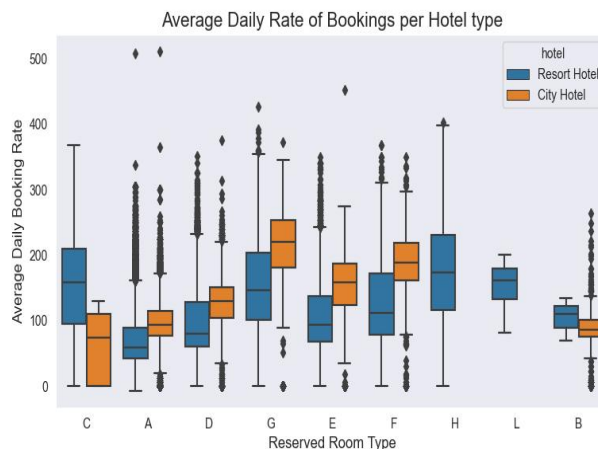
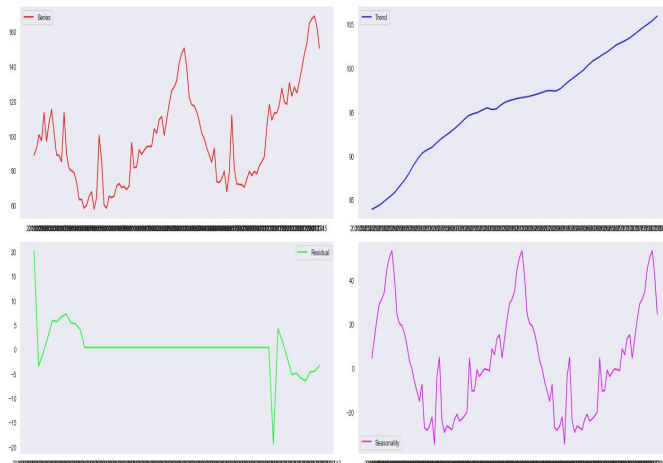


Fig. 1. Average daily rate of bookings per hotel type

At this stage, the exploratory data analysis is performed in order to check through data features and find whether there are any outliers. A column "FullDate" is built to facilitate temporal analysis. It is formed by joining the arrival date year and week number together to indicate Week 43 of the year 2015, for instance, by "201543". Next, the ADR is averaged and after being made into a time series then, the object is created every week. To separate the joined data into fractions such as training and testing sets, a data frame is produced [8]. The dataset is split into two sets for training and testing: the training set is made up of rows 1-100, and the testing set is made up of the final rows 101 to 125. The main objective of this preprocessing stage is to produce the maximum possible demand for real-time hotel bookings thereby forming a compelling groundwork for further investigation in this regard..

### 4.3 Data decomposition

Data decomposition is done to figure out the presence of trend and seasonality components in the dataset. Trend tells about the ultimate trend of the data, showing it whether it is constantly going up, falling, or staying the same over time. On the contrary, several factors influence seasonality which constantly repeat the pattern at regular intervals as well as they are impacted by factors such as seasons or festivals [12]. The decomposition certainly does show these trends in the process of close consideration — a verdict that seems to be ineluctable. Since a residual chart assesses the rest of the decomposition in compliance with the trend and seasonality of the given period, in addition to these two basic properties, we can conclude that the decomposition is correct. Such data understanding of the taking times of demands for hotel bookings is crucial for analysing the temporal dynamics of these demands, which makes it possible to go further with optimisation efforts.



**Fig. 2.** Data decomposition

### 4.4 Model selection

Model selection is done after analysing the train plot produced with a moving average model (ARIMA) using Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF) plots to ascertain fit [4]. Afterwards, augmented Dickey-Fuller (ADF) test is administered to tackle the non-stationarity of the sensitivity data; thereafter, deviations are pinpointed and taken out. The data is then differenced with the exogenous factors, and SARIMAX [6] model (temporal Autoregressive Integrated Moving Average with Exogenous) is eventually tailored. On top of that, a choice of a seasonality factor brings a trendline deletion that already exists in data. The sign of the data is used to identify the stationary behaviour of the variables. Finally, the SARIMAX model is shown to be suitable for real-time forecasting of hotel booking demand via the validation of the augmented Dickey-Fuller test together with its suitability in a stationary series.

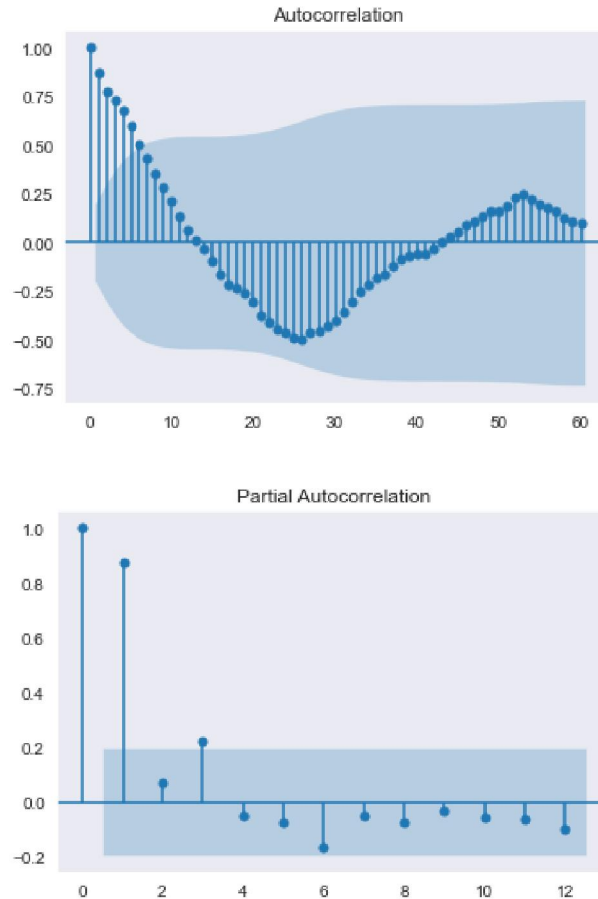


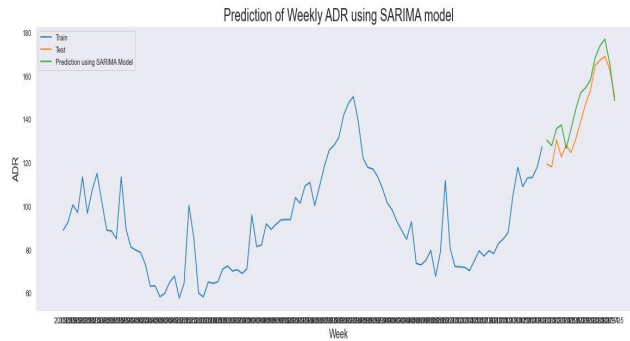
Fig. 3. ACF and PACF plots

#### 4.5 Data modelling

To find the perfect model for the dataset, the Auto-Arima function is implemented, which picks ARIMA(0,1,1) (0,1,0) and the lowest Akaike Information Criterion (AIC) score. Besides, the SARIMAX function or the corresponding model function is used. The final result of SARIMAX will be submitted to the Ljung-Box Test which is for checking their residuals auto-correlation [11]. We plot, where actually observed values compared with the predicted using a daily rate are displayed for weekly data to assess model's performance. The outcome shows a RMSE of 8.75 that is a quantifier form of accuracy for the model. The comprehensive data modelling nature operates on revealing an idea for an optimal strategy which targets the real-time hotel booking demand.

SARIMAX Results					
Dep. Variable:	y			No. Observations:	100
Model:	SARIMAX(0, 1, 1)x(0, 1, [], 52)			Log Likelihood	-182.876
Date:	Sun, 24 Apr 2022			AIC	369.751
Time:	00:58:06			BIC	373.452
Sample:	0			HQIC	371.144
					- 100
Covariance Type:	opg				
	coef	std err	z	P> z	[0.025 0.975]
ma.L1	-0.5853	0.098	-5.967	0.000	-0.778 -0.393
sigma2	139.0961	23.321	5.964	0.000	93.387 184.805
Ljung-Box (L1) (Q):	0.30	Jarque-Bera (JB):	5.74		
Prob(Q):	0.58	Prob(JB):	0.06		
Heteroskedasticity (H):	0.55	Skew:	-0.60		
Prob(H) (two-sided):	0.25	Kurtosis:	4.22		

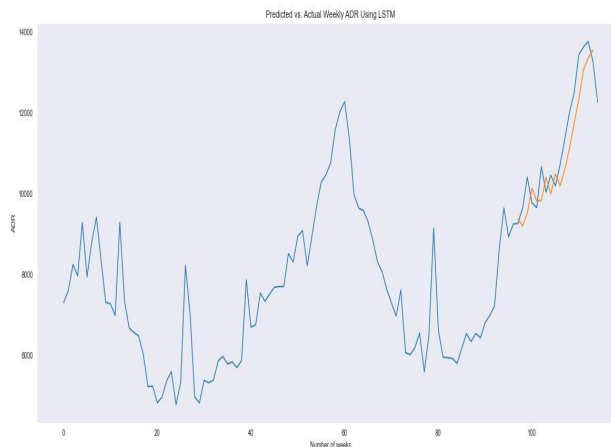
Fig. 4. Model summary



**Fig. 5.** ADR plot using SARIMA model

#### 4.6 Prediction using LSTM

The Transformation of a data frame into an array is a part of the process and mainly done by the LSTM method. Finally, the dataset is partitioned into training and test subsets with the test subset comprising 20% of the data and 80% going to training class. The MinMaxScaler is to be incorporated in the standardisation of the data, and a lookback period of 5 is to be employed. Next, the Adam optimizer is utilised and mean square error is set as the loss function to train the LSTM model for 100 iterations. To introduce to our model the size of the dataset we'll visualise the train loss and validation loss. And the corresponding results which are the root mean square errors (RMSE) that represents the difference between the actual daily rates and the projected ones shows a very low value of 8.17 which is actually smaller than the former version. Such LSTM application shows possibilities how accurately it may predict the filing of hotel booking demands in real time and thus has a place in optimising the operations [12].



**Fig. 6.** Predicted value vs ADR using LSTM

### V. BINARY CLASSIFICATION

The main aim of the binary classifications task is to predict whether the hotel booking will be cancelled or not. Using a machine learning algorithm, this predictive modelling approach divides booking instances into two classes: It has the ability to pass two modes, therefore mark them as cancelled or not cancelled class [16]. By testing the different parameters of reservation lead time, arrival date and guest characteristics and getting the machine to learn what settings are most likely to point to trip cancellations we are applying a reinforcement classification method.

#### 5.1 Feature Engineering

In addition to the existing features, fifteen new features added to the dataset with the help of feature engineering approach. For the ease of identification, a prefix "new\_" is attached to the dataset for distinction of these attributes. For instance, "new\_is\_weekend" takes 1 for the weekend when the booking date is not weekdays and 0 for weekdays.

instance, "new\_is\_family" would be displayed to highlight the fact that there are mothers or small children in the reservation [9]. This procedure produces 44 features all in all which enhances the content richness of the data set and provides a more comprehensive and deep analysis to finalise demand for customers who would need a hotel booking instantaneously.

hotel	object	customer_type	object
is_canceled	int64	adr	float64
lead_time	int64	required_car_parking_spaces	int64
arrival_date_year	int64	total_of_special_requests	int64
arrival_date_month	object	reservation_status_date	int64
arrival_date_week_number	object	new_is_family	int64
arrival_date_day_of_month	int64	new_room_difference	int64
stays_in_weekend_nights	int64	new_total_people	float64
stays_in_week_nights	int64	new_total_stay_day	int64
adults	float64	new_month	int64
children	int64	new_arrival_date	int64
babies	object	new_PMS_entering_date	int64
meal	object	new_special_req_status	int64
country	object	new_dist_channel_type	object
market_segment	object	new_room_difference_cat	float64
distribution_channel	int64	new_is_weekend	int64
is_repeated_guest	int64	new_is_weekday	int64
previous_cancellations	int64	new_is_weekend_and_weekdays	int64
previous_bookings_not_canceled	int64	new_want_parking_space	int64
reserved_room_type	int64	new_adr_per_person	float64
assigned_room_type	object	dtype: object	
booking_changes	int64		
deposit_type	int64		
days_in_waiting_list	int64		

Fig. 7. Feature engineering

**5.2 Data encoding**

The columns having a lot of classes like the ones of 300-class level, are difficult to analyse and this is resolved by data encoding. Among other approaches for simplifying representation, object-type columns containing binary values are undergoing label encoding. One-hot encoded binary columns for any class in columns that have three to twelve classes, which enhances categorical representation, is used [5]. Also, to maintain unity in scale among features, standardisation is applied to the numeric columns by means of the Standard Scaler. Through such methods of encoding real-time demand for hotel booking effectively can be made possible and also data processing and analysis will be more efficient.

**5.3 Model selection for binary classification**

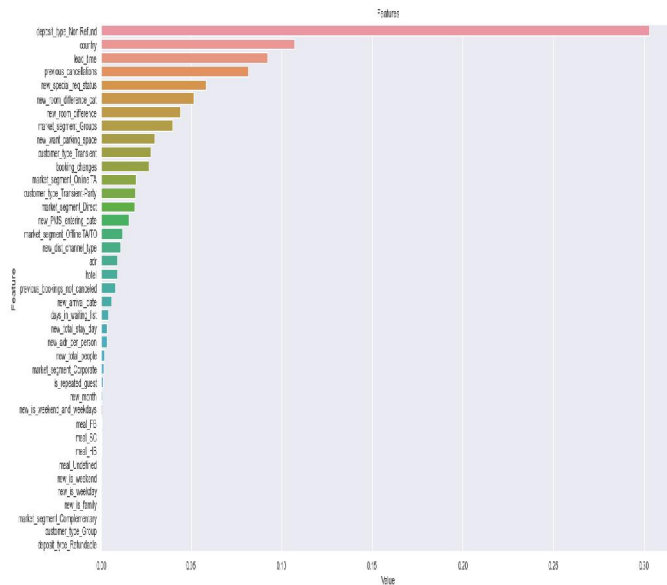


Fig. 8. Feature plot

In the model selection approach for binary classification six different models are considered and their performance metrics including F1 score, accuracy, and ROC & AUROC scores, are evaluated. The obtained CART model is deep to

1 and with min division to 2 reaches F1 score 0.4775 and accuracy 74.98%. KNN, however, involved 50 neighbours and had the F1 score of 0.4783 and the accuracy of 62.08%. At maximum depth 5, maximum features 5, and 100 estimators, Random Forest hits an F1 score 0.4945, it has 74.82% accuracy, and has an outstanding ROC AUC 0.8346. AdaBoost with 10 estimators and a learning rate of 0.98 yields an F1 of 0.5875 and 75.51% accuracy. ROC & AUC is 0.8450. In the same way, a gradient boosting machine with 100 estimators and learning rate of 0.01, provides 0.4882 F1 Score, 75.25% accuracy, and 0.8504 for the ROC AUC. With the last experiment, the Light gradient boosting classifier algorithm based on gradient boosting machine with 100 estimators, a learning rate of 0.01, and the column sample by tree parameter of 0.7 will show an F1 score of 0.5607 and an accuracy of 69.86% which has the ROC AUC of 0.7895 Accordingly, the gradient boosting machine model came out as the best model under the binary classification task because of its high accuracy as well as high ROC and AUC scores.

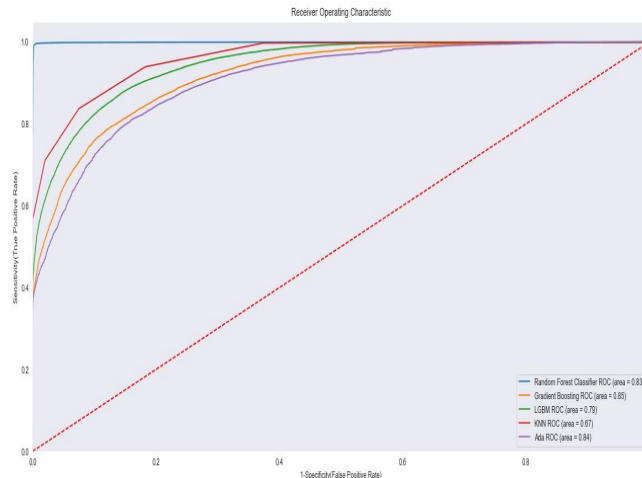


Fig. 9. ROC plot

## VI. CONCLUSION

In conclusion, LSTM outperformed SARIMA model in terms of the forecasting accuracy for the average daily rate since LSTM is more sensitive to sequence length of various temporal structures. The ADABOOST method, however, recorded the highest performance levels for the six classification models indicating that it is strong for classification tasks [12]. With the healthy prediction model of 1,80,000 bookings in the test data, a good estimation of the number of cancellations, i.e. 12,000 were made for better business decisions. Through the use of the LSTM model, a specific daily rate was obtained for such bookings which in turn makes a clear revenue projections forecast possible. Considering the limitations of this model, now it is important to make it more versatile and flexible. This versatility can enable it to be applied to larger datasets of different accommodation types like hostels, Airbnbs, and so on [17]. This forecasting ability is helpful in eliminating uncertainties for purposeful long-term planning and operational potentials of an industry that is continuously changing.

## VII. FUTURE WORK

In the future, the model can be developed to incorporate the new set features of many different accommodation types like hostels, Airbnbs, Vrbo, etc., in different locations. This upgraded model would yield weekly average daily rates forecasts that will be more adjustable to a wider variety of lodging options and consequently will go beyond the limits of a hospitality industry [3]. As an auxiliary project, had the LSTM exhibited a superior forecast execution as compared to SARIMA for predicting the average daily rate and AdaBoost given a successful application to binary classification, then the subsequent optimization and fine-tuning of the models could be explored. The advent of intelligent machine learning and deep learning technologies opens up possibilities of using them as reinforcing factors in improving the accuracy and anchorage of prediction models for hotel booking demand optimization. Moreover, other upcoming development projects can target overseeing numerous varied datasets belonging to multiple entities to come up with



inventive models that can accurately diagnose hotel reservation cancellations. Besides that, considering ensemble learning models- or hybrid ones could bring in some extra boost in precision and generalisation tasks- performance. This powerful strategy would be incorporated and continue to pave the way for predictive analytics in hospitality, making decision making and strategy formulation easy to ensure maximum output.

#### REFERENCES

- [1] N. Phumchusri and P. Ungtrakul, "Hotel daily demand forecasting for high-frequency and complex seasonality data: a case study in Thailand," *Tourism Economics*, vol. 26, no. 7, pp. 1080-1100, Dec.2020. DOI:<https://doi.org/10.1057/s41272-019-00221-6>.
- [2] Nuno Antonio, Ana de Almeida, and Luis Nunes. "Hotel booking demand datasets". In: *Data in brief* 22 (2019), pp. 41–49.
- [3] Aalen, P., Iversen, E. K., & Jakobsen, E. W. (2019). Exchange rate fluctuations and demand for hotel accommodation: Panel data evidence from Norway. *Scandinavian Journal of Hospitality and Tourism*, 19(2), 210–225.
- [4] Choi, J.-G. (2007). Developing a restaurant industry business cycle model and analysing industry turning points. *Journal of Global Business and Technology*, 3(1), 40–48.
- [5] Enz, C. A., Canina, L., & Lomanno, M. (2009). Competitive pricing decisions in uncertain times. *Cornell Hospitality Quarterly*, 50(3), 325–341.
- [6] Gallagher, M., & Mansour, A. (2000). An analysis of hotel real estate market dynamics. *Journal of Real Estate Research*, 19(2), 133–164.
- [7] Nuno Antonio, Ana De Almeida, and Luis Nunes. "Predicting hotel booking cancellations to decrease uncertainty and increase revenue". In: *Tourism & Management Studies* 13.2 (2017), pp. 25–39.
- [8] Parul Gupta, "Analysis of Customer Satisfaction of the Hotel Industry in India Using Kano Model & QFD," *International Journal of Research in Commerce, IT & Management*, vol. 2, no. 1, pp. 1-10, January 2012, ISSN 2231-5756.
- [9] Grebler, L., & Burns, L. S. (1982). Construction cycles in the United States since World War II. *Real Estate Economics*, 10(2), 123–151.
- [10] Manoranjan Dash and Huan Liu. "Feature selection for classification of hotel industry", *Intelligent data analysis* 1.3 (1997), pp. 131–156.
- [11] Han-Chen Huang, Allen Y Chang, Chih-Chung Ho, et al. "Using artificial neural networks to establish a customer-cancellation prediction model". In: *Przeegląd Elektrotechniczny* 89.1b (2013), pp. 178–180.
- [12] Dolores Romero Morales and Jingbo Wang. "Forecasting cancellation rates for services booking revenue management using data mining". In: *European Journal of Operational Research* 202.2 (2010), pp. 554–562.
- [13] Sheather, S. (2019). *A modern approach to regression with R for forecasting hotel booking Cancellations*. Springer Science & Business Media.
- [14] Breffni M Noone and Chung Hun Lee. "Hotel overbooking: The effect of overcompensation on customers' reactions to denied service". In: *Journal of Hospitality & Tourism Research* 35.3 (2011), pp. 334–357.
- [15] Song, H., & Li, G. (2008). Tourism demand modelling and forecasting—A review of recent research. *Tourism Management*, 29(2), 203–220.
- [16] Wallis, G. (2019, June). Meet the money: Sunny skies continue in Hotel Business, 28(8), 20.
- [17] M. Chattopadhyay and S. K. Mitra, "Determinants of revenue per available room: Influential roles of average daily rate demand seasonality and yearly trend", *Int. J. Hospitality Manage.*, vol. 77, pp. 573-582, Jan. 2019.