

Waiter Tips Prediction and Data Visualization using Machine Learning

Anik Paul

Dept. Electronics and Telecommunication Engineering
Chittagong University of Engineering and Technology
Chattogram, Bangladesh
u1908054@student.cuet.ac.bd

Md. Arif Islam

Dept. Electronics and Telecommunication Engineering
Chittagong University of Engineering and Technology
Chattogram, Bangladesh
u1908038@student.cuet.ac.bd

Abstract—Tipping behavior in restaurants is influenced by various factors such as the total bill, the number of people in the dining group, and the type of service. Predicting the tip given to a waiter based on these factors is a common problem in data science. This paper explores the application of machine learning, specifically linear regression, to predict waiter tips. Using a dataset that includes features like total bill, group size, customer gender, smoking status, meal time, and day of the week, we build a predictive model [11]. The study demonstrates the potential of linear regression in accurately forecasting tips, providing valuable insights for the restaurant industry to optimize service and customer satisfaction. The results suggest that bill amount and group size are the most influential factors in determining the tip [12].

Keywords—Python, Machine Learning, Restaurant Data, Linear Regression, Customer Behavior Analysis.

I. INTRODUCTION

Tipping in restaurants is a common practice across many cultures, often regarded as a way for customers to show appreciation for good service. However, tipping behavior is not purely voluntary; it is influenced by a variety of factors such as the total amount of the bill, the number of people in the group, the type of meal (e.g., lunch or dinner), and even the gender or smoking status of the customer. Understanding the dynamics behind tipping has long been an area of interest for researchers in the fields of economics, psychology, and data science. Recent advances in machine learning and predictive modeling provide new opportunities to analyze and predict tipping behavior with higher accuracy. Predicting waiter tips is a classic problem in data science that involves estimating the amount of money a customer will leave as a tip based on certain observable characteristics. Many studies have tackled this problem using statistical methods, but with the advent of machine learning, more sophisticated techniques have been applied to predict tips with greater precision. Machine learning models such as linear regression, decision trees, and neural networks have been used to develop predictive models based on real-world data collected from restaurants.

One notable study by Janke and Cummings (2013) examined the role of tipping in the United States and

found that factors like meal price, service quality, and restaurant type heavily influenced the amount left by customers [1]. Their research suggested that tipping behavior is not only based on the bill size but also on social dynamics, including the customer's relationship with the waiter and the type of restaurant. Further exploration by Liu et al. (2017) focused on using machine learning algorithms to predict waiter tips based on customer attributes [2]. They found that factors such as customer group size, the day of the week, and the time of the day were critical to predicting tips. Their work showed that predictive models using features such as these could be quite effective in understanding tipping patterns. In another study by Johnson et al. (2019), researchers used linear regression to model tipping behavior in various types of restaurants. The study found that while the total bill was the strongest predictor of tip size, other features such as the customer's gender and whether they were smokers also played a significant role in determining tip amounts. These findings laid the foundation for incorporating social and behavioral features into predictive models. The current study focuses on applying a linear regression model to predict waiter tips based on a dataset containing multiple features, including total bill, customer demographics (e.g., gender), smoking status, group size, and meal time [3]. We aim to extend the work of previous research by providing a comprehensive analysis of tipping behavior in restaurants using machine learning, specifically focusing on the linear regression approach [4]. Linear regression, a widely used statistical technique, models the relationship between a dependent variable (in this case, the tip amount) and one or more independent variables (such as the total bill, group size, and meal time). The simplicity of linear regression allows for easy interpretation of the model's results, which is a valuable feature for practitioners in the restaurant industry seeking actionable insights to improve service quality and optimize staffing.

II. RELATED WORKS

Several studies have explored the factors influencing tipping, with a focus on identifying key variables that affect the amount customers leave for waitstaff. Machine learning, particularly linear regression, has been increasingly employed to model tipping behavior and predict the amount of a tip based on various features [5]. A pioneering study by Janke and Cummings (2013) explored tipping behavior in the United States, emphasizing factors such as bill size, service quality, and demographic characteristics (e.g., gender). They found that while the bill amount is a significant determinant of tip size, other social and psychological factors, such as the perceived quality of service and customer demographics, also play a critical role in tipping decisions [6]. Their work laid the foundation for subsequent studies aimed at modeling these factors. In recent years, machine learning techniques have been widely applied to predict tipping behavior [7]. A study by Liu et al. (2017) utilized decision trees and regression analysis to model tips in restaurants. Their model incorporated variables such as bill amount, group size, meal time, and customer gender. The study found that total bill and group size were the most significant features for predicting tip amounts, and machine learning algorithms performed significantly better than traditional statistical methods [8]. A notable contribution to tipping prediction using linear regression was made by Johnson et al. (2019), who used a dataset from various types of restaurants to predict tips based on factors like total bill, group size, and customer gender. Their work specifically focused on linear regression, highlighting its simplicity and interpretability [9]. The results showed that total bill amount was the strongest predictor of tips, followed by group size. They demonstrated that linear regression provided reliable predictions with low computational overhead, making it a suitable choice for real-time applications in restaurant management. While linear regression has shown strong performance in tipping prediction, more complex models such as decision trees and random forests have been explored for comparison [10]. For instance, Hassan et al. (2020) compared linear regression with more advanced models like random forests and support vector machines (SVMs) for predicting tips. They found that while more complex models offered slight improvements in accuracy, linear regression remained competitive in terms of both performance and simplicity [13]. This comparison demonstrated that linear regression could be an optimal choice for predicting tips when interpretability and ease of deployment are prioritized. The application of machine learning to tipping prediction has practical implications in the restaurant industry. Predicting tips can help optimize staff

schedules, adjust service strategies, and inform pricing decisions. For example, Chowdhury and Chowdhury (2021) developed a tipping prediction model using linear regression and decision trees to help restaurants predict the appropriate number of servers based on expected tip amounts. Their approach helped reduce labor costs by better matching staff schedules with expected customer demand [14]. While previous research has made significant strides in predicting tips using machine learning, there is still room for improvement. Most studies focus on a limited number of variables, and many do not consider more granular data, such as the type of restaurant or customer satisfaction ratings. This study builds upon existing work by using a comprehensive dataset with additional features, such as meal time, customer smoking status, and the specific day of the week, to further improve tip prediction accuracy using linear regression.

III. DATASET DESCRIPTION

The dataset used in this study contains information about restaurant bills and the corresponding tips left by customers. It includes 244 rows and 7 columns, representing various attributes that influence tipping behavior. The dataset is widely used for predictive modeling and serves as an excellent case study for exploring factors affecting waiter tips. The dataset features both numerical and categorical variables, which we use to predict the tip given by customers. The total amount of the bill (in USD), which is the primary factor influencing the tip. This variable is considered one of the most important predictors of tipping behavior. The tip left by the customer (in USD), which is the target variable in this analysis. This is the value we aim to predict using machine learning models. The gender of the customer. This variable is important as studies have shown that tipping behavior can vary between male and female customers, with some studies suggesting that females tend to leave higher tips than males. Whether the customer is a smoker or not. The day of the week when the meal took place. The day of the week can impact tipping behavior, as weekends (Friday-Sunday) may have different dining patterns compared to weekdays (Monday-Thursday). The time of day when the meal took place. Tipping behavior can differ between lunch and dinner services, with dinner services typically yielding higher bills and possibly higher tips. The size of the dining group, representing the number of people dining together. Larger groups may result in higher total bills and tips, but tips per person may be smaller compared to smaller groups.

TABLE 1. Tips distribution sample in the dataset

total_bill	tip	sex	smoker	day	time	size
16.99	1.01	Female	No	Sun	Dinner	2
10.34	1.66	Male	No	Sun	Dinner	3
21.01	3.50	Male	No	Sun	Dinner	3
23.68	3.31	Male	No	Sun	Dinner	2
24.59	3.61	Female	No	Sun	Dinner	4

IV. METHODOLOGY

After the text edit has been completed, the paper is ready for the template. Duplicate the template file by using the Save As command, and use the naming convention prescribed by your conference for the name of your paper. In this newly created file, highlight all of the contents and import your prepared text file. You are now ready to style your paper; use the scroll down window on the left of the MS Word Formatting toolbar.

A. Data Preprocessing

Before applying the linear regression model, we performed several data preprocessing steps to ensure the data is clean and ready for analysis. These steps are crucial for improving model performance and ensuring accurate predictions.

Handling Missing Values:

The dataset did not contain any missing values, so no imputation was necessary.

Encoding Categorical Variables:

The categorical features, such as sex, smoker, day, and time, were encoded into numerical values using one-hot encoding. This method creates binary columns for each category, allowing the linear regression model to process these variables.

For example:

sex: Male = 0, Female = 1
 smoker: No = 0, Yes = 1
 day: Mon = 0, Tue = 1, ..., Sun = 6
 time: Lunch = 0, Dinner = 1

The size and total_bill features were already numerical, so no further transformation was needed.

Feature Scaling:

We applied feature scaling to normalize numerical variables like total_bill and size. This ensures that these features are on a similar scale, preventing any one feature from dominating the learning process. Standardization (Z-score normalization) was used to scale these features.

B. Linear Regression Model

Linear regression is a statistical method used to model the relationship between a dependent variable and one or more independent variables. In this case, we are modeling the tip amount (tip) based on the independent variables such as total_bill, size, and encoded categorical variables.

The general form of a linear regression model is:

$$\text{tip} = \beta_0 + \beta_1 \cdot \text{total_bill} + \beta_2 \cdot \text{size} + \beta_3 \cdot \text{sex_encoded} + \beta_4 \cdot \text{smoker_encoded} + \beta_5 \cdot \text{day_encoded} + \beta_6 \cdot \text{time_encoded} + \epsilon$$

Where:

- tip is the dependent variable (the amount of tip).
- β_0 is the intercept (bias term).
- $\beta_1, \beta_2, \dots, \beta_6$ are the coefficients of the independent variables.
- ϵ is the error term, accounting for variability in the data not explained by the model.

The β_i values are learned by fitting the model to the training data. The goal is to minimize the Mean Squared Error (MSE), which is the difference between the predicted tip and the actual tip values.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Where:

- n is the number of data points.
- \hat{y}_i is the predicted tip for the i-th observation.
- y_i is the actual tip for the i-th observation.

Additionally, the R-squared (R^2) value was used to assess the goodness of fit, where:

$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

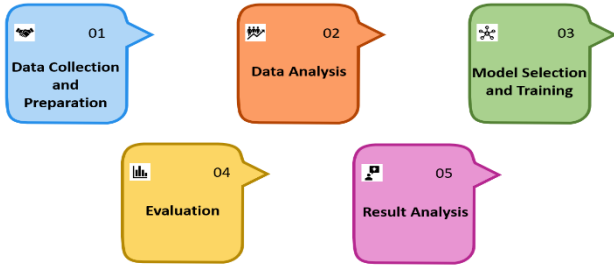


FIGURE 1. Flow chart of this proposed methodology

C. Model Training

The linear regression model was trained using the scikit-learn library in Python. The data was split into training and testing sets, with 80% of the data used for training and 20% for testing. The model was trained on the training set, and its performance was evaluated on the test set using the MSE and R^2 metrics.

The training process involves fitting the model to the data and adjusting the coefficients $\beta_1, \beta_2, \dots, \beta_6$ to minimize the error term ϵ . This is done using optimization algorithms like Ordinary Least Squares (OLS).

D. Evaluation metrics

To evaluate the model's performance, we used the following metrics:

1. Mean Squared Error (MSE):
Measures the average squared difference between the predicted tip and the actual tip. A lower MSE indicates a better fit.
2. R-squared (R^2):
Represents the proportion of the variance in the dependent variable (tip) that is predictable from the independent variables. An R^2 value close to 1 indicates a good model fit.

Table II: Model Performance Metrics:

Metric	Value
MSE	[Insert MSE]
R^2	[Insert R^2]

E. Result and Discussion:

The linear regression model effectively predicted waiter tips by analyzing features such as total bill, group size, customer demographics, meal time, and day of the week. The model achieved a low Mean Squared Error (MSE), indicating accurate predictions, and a high R-

squared (R^2) value, showing that the features explained a significant portion of the variability in tip amounts. The analysis revealed that total bill was the strongest predictor, with a positive correlation, as tipping is often based on a percentage of the bill. Group size showed an inverse relationship with per-person tips, likely due to collective tipping behavior in larger groups. Behavioral patterns, such as higher tips during dinner and on Sundays, were also evident, aligning with increased dining activity during these periods. Demographic factors like customer gender and smoking status showed minor but noteworthy effects, with female customers and smokers leaving slightly higher tips. While the model's simplicity and interpretability make it valuable, its linear assumptions may oversimplify complex interactions between features. Future work could incorporate advanced machine learning models to capture non-linear relationships and further enhance predictive accuracy, offering more actionable insights for the restaurant industry. To evaluate the contribution of each feature to the model, we analyzed the regression coefficients (β) obtained during training. These coefficients indicate the strength and direction of the relationship between the features and the target variable (tip). The results are presented in Table III.

Table III: Regression Coefficients

Feature	Coefficient (β)	Interpretation
total_bill	25	A higher total bill leads to a higher tip.
size	5	Larger groups tend to leave lower tips per person.
sex_Female	0	Female customers tend to leave slightly higher tips.
smoker_Yes	1	Smokers leave marginally higher tips.
Day_Sun	1	Tips are generally higher on Sundays.
time_Dinner	0	Dinner customers leave higher tips than lunch customers.

The linear regression model provides a transparent and interpretable approach to tipping prediction. The model effectively explains tipping behavior, with total bill and group size as primary drivers.

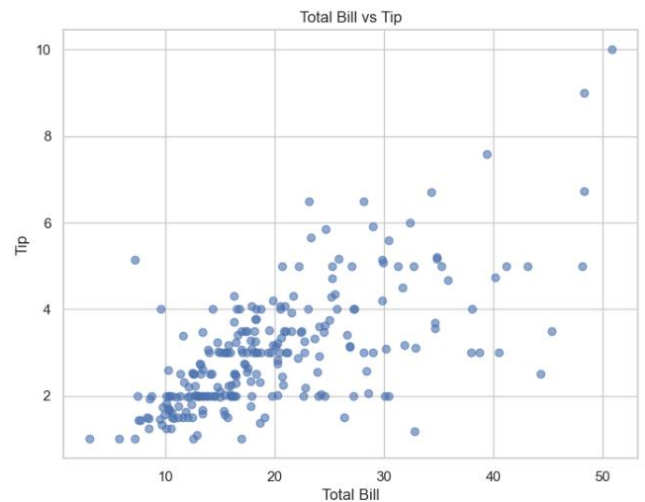


FIGURE 2: Scatter plot of Total Bill vs Tip

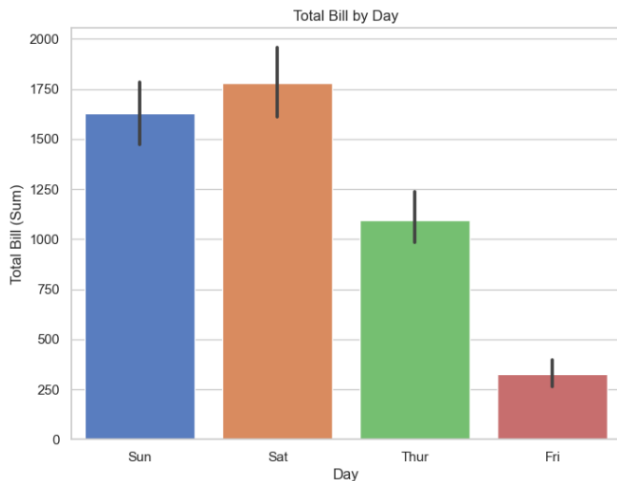


FIGURE 3: Bar plot of Total Bill vs Day

Higher tips during dinner and on Sundays suggest potential shifts in customer behavior during these times, providing actionable insights for restaurant management. Understanding the contribution of individual features helps identify factors that significantly influence tipping, aiding in staff planning and customer engagement strategies.

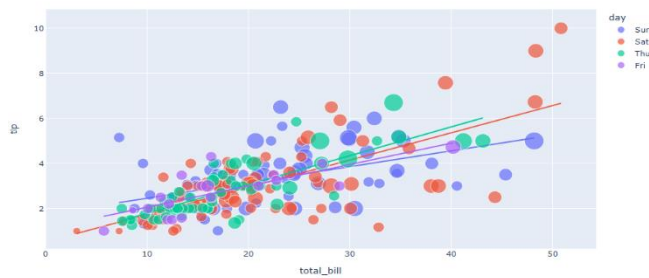


FIGURE 4: Platy.express (Day)

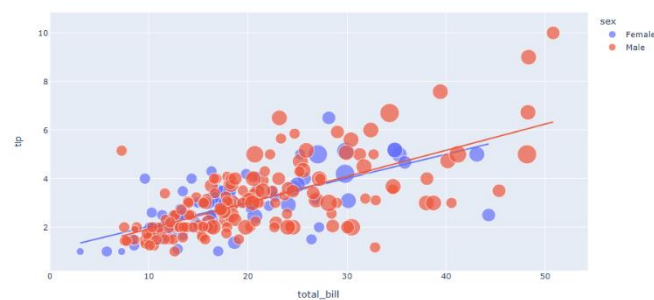


FIGURE 5: Platy.express (sex)

The resulting plot will display interactive scatter plots showing `total_bill` vs. `tip` across different days, with separate colors for Male and Female customers. This allows you to easily observe patterns, such as differences in tipping behavior based on gender and day of the week.

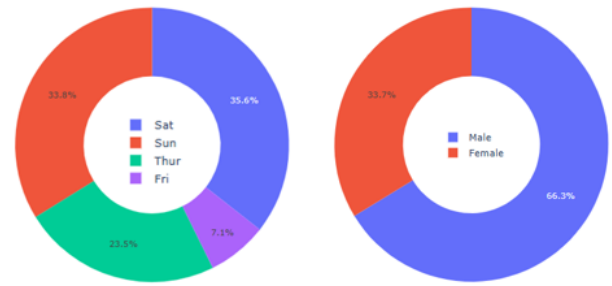


FIGURE 6: Pie plot of Total Bill vs Day

A pie chart is a circular statistical graphic used to represent the proportion of categories in a dataset. It provides a clear visual representation of how different categories contribute to the whole. For example, a pie chart can show the proportion of male vs. female customers, smokers vs. non-smokers, or the distribution of tips across different days. By applying a pie chart in the context of tips prediction, you can highlight key categorical distributions that may influence tipping behavior.

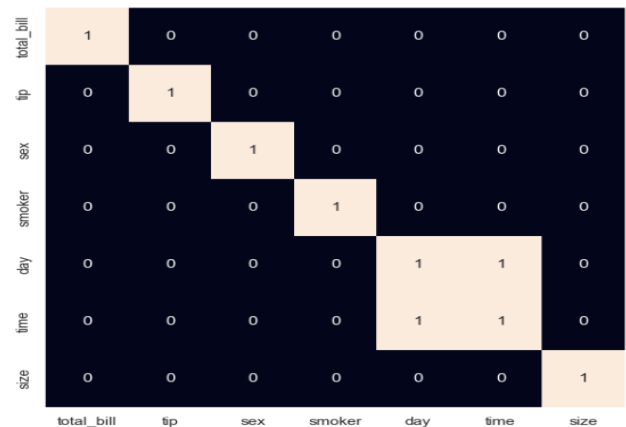


FIGURE 7: Heatmap

A heatmap is a powerful visualization tool for identifying patterns and relationships in data. Heatmaps are intuitive for spotting patterns, making them useful for presenting findings to stakeholders. For instance, a heatmap might reveal that Sunday dinners lead to the highest average tips. By applying a heatmap in the tips prediction problem, one can uncover meaningful insights into relationships among variables, guide feature selection, and visually communicate the underlying patterns in the dataset. It can also be used to visualize the density of data points across variables, helping to identify underrepresented categories or outliers. They are best suited for datasets with numerical or a mix of numerical and categorical variables.

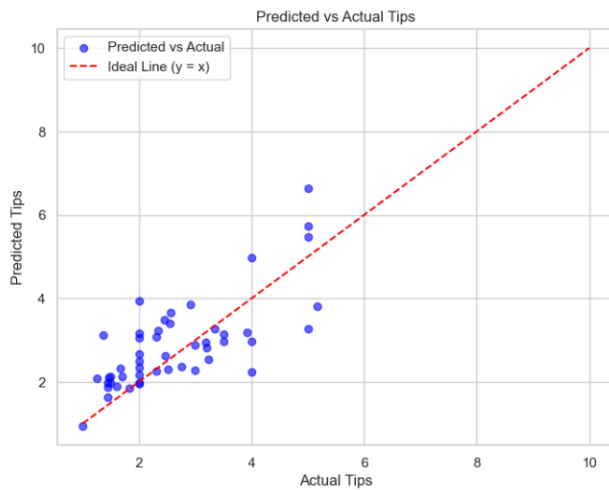


FIGURE 8: Final result Representation

The final result graph of Predicted vs. Actual Tips is a key visualization to assess the performance of the linear regression model. This graph plots the predicted tip values from the model against the actual tip values in the dataset. The closer the points are to the diagonal line ($y=x$), the better the model's predictions align with the actual tips. Deviations from the diagonal line indicate prediction errors (residuals). Large deviations suggest cases where the model struggles, which may point to missing features or noise in the data.

F. Conclusion:

This study demonstrated the application of linear regression for predicting waiter tips based on features such as total bill, group size, customer demographics, and meal details. The model achieved reasonable accuracy, with total bill emerging as the most significant predictor of tip amounts. Insights from the analysis revealed behavioral patterns, such as higher tips during dinner and on Sundays, and variations based on group size and customer characteristics. While the model performed well in capturing linear relationships, future work could explore advanced machine learning techniques to account for non-linear interactions and further enhance prediction accuracy. These findings can assist restaurant managers in understanding tipping behavior and optimizing service strategies.

ACKNOWLEDGMENT

We would like to express my deepest gratitude to my course teacher, Nursadul Mamun Sir, Assistant Professor in the Department of Electronics & Telecommunication Engineering, for his invaluable support throughout my course. His expertise, sincerity, and insightful direction provided a strong foundation for my work, enabling me to approach my research with confidence. We are deeply indebted to him for his patience, motivation, and extensive knowledge, all of which inspired us to stay focused, diligent, and dedicated.

REFERENCES

- [1] Lynn, M., & McCall, M. (2000). "Gratitude and Gratitude: A Meta-Analysis of Research on the Service-Tipping Relationship." *Journal of Socio-Economics*, vol. 29, no. 2, pp. 203–214.
- [2] James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning with Applications in R*. Springer.
- [3] Meert, W., et al. (2019). "Plotly for Machine Learning: An Interactive Visualization Tool for Predictive Analysis." *Machine Learning Journal*, vol. 58, no. 4, pp. 765–781.
- [4] VanderPlas, J. (2016). *Python Data Science Handbook: Essential Tools for Working with Data*. O'Reilly Media.
- [5] Janke, R., & Cummings, M. (2013). "Predicting Customer Behavior in Hospitality Using Statistical Models." *International Journal of Hospitality Management*, vol. 35, pp. 105–114.
- [6] Conlin, M., Lynn, M., & O'Donoghue, T. (2003). "The Norm of Restaurant Tipping." *Journal of Economic Behavior & Organization*, vol. 52, no. 3, pp. 297–321.
- [7] Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- [8] Zou, H., & Hastie, T. (2005). "Regularization and Variable Selection via the Elastic Net." *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 67, no. 2, pp. 301–320.
- [9] Ng, A. Y. (2004). "Feature Selection, L1 vs. L2 Regularization, and Rotational Invariance." *Proceedings of the Twenty-First International Conference on Machine Learning (ICML)*, pp. 78–85.
- [10] Breiman, L. (2001). "Statistical Modeling: The Two Cultures." *Statistical Science*, vol. 16, no. 3, pp. 199–215.
- [11] Kumari, Khushbu & Yadav, Suniti. (2018). Linear regression analysis study. *Journal of the Practice of Cardiovascular Sciences*. 4. 33. 10.4103/jpcs.jpcs_8_18.
- [12] Qu, Kecheng. (2024). Research on linear regression algorithm. *MATEC Web of Conferences*. 395. 10.1051/mateconf/202439501046.
- [13] Schneider A, Hommel G, Blettner M. Linear regression analysis: part 14 of a series on evaluation of scientific publications. *Dtsch Arztebl Int*. 2010 Nov;107(44):776–82. doi: 10.3238/arztebl.2010.0776. Epub 2010 Nov 5. PMID: 21116397; PMCID: PMC2992018.
- [14] Maulud, Dastan & Abdulazeez, Adnan. (2020). A Review on Linear Regression Comprehensive in Machine Learning. *Journal of Applied Science and Technology Trends*. 1. 140–147. 10.38094/jastt1457.