



Applied Statistics

Final Project

Miguel Bourgin Gonçalves
nº202107127

Index

1	Introduction	3
1.1	Linear Regression	3
1.2	Logistic Regression	3
2	Methodology	4
3	Model Evaluation and Interpretation	4
4	Data Pre-Processing	5
5	Exploratory Data Analysis	6
6	Model Selection Process	8
7	Analysis of Variables	10
8	Conclusion	12

Abstract

This project aims to apply statistical methods to analyze specific data and extract relevant information. A quantitative approach will be used, focusing on descriptive and inferential analysis. The study will consider "tip" as the response variable, examining its relationship with various predictors. It will cover the methodology used for data analysis and model evaluation, followed by the presentation of the results obtained, an interpretation and discussion of the results, and the final conclusions.

1 Introduction

1.1 Linear Regression

Linear regression is a method used to model the relationship between a continuous response variable Y and one or more predictor variables X . It assumes a linear relationship, where the response is predicted as a weighted sum of the predictors, with the goal of minimizing the difference between the predicted and actual values of Y .

Simple Linear Regression

Simple linear regression lives up to its name: it is a very straightforward approach for predicting a quantitative response Y on the basis of a single predictor variable X . It assumes that there is approximately a linear relationship between X and Y . Mathematically, we can write this linear relationship as

$$Y \approx \beta_0 + \beta_1 X. \quad (1.1)$$

In equation (1.1), β_0 and β_1 are two unknown constants that represent the intercept and slope terms in the linear model. Together, β_0 and β_1 are known as the *model coefficients* or *parameters*. Once we have used our training data to produce estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ for the model coefficients, we can predict future values of Y on the basis of a particular value of X by computing

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x, \quad (1.2)$$

where \hat{y} indicates a prediction of Y on the basis of $X = x$.

Multiple Linear Regression

Simple linear regression is a useful approach to predicting a response on the basis of a single predictor variable. However, in practice, we often have more than one predictor. Instead of fitting a separate simple linear regression model for each predictor, a better approach is to extend the simple linear regression model so that it can directly accommodate multiple predictors. We can do this by giving each predictor a separate slope coefficient in a single model.

In general, suppose that we have p distinct predictors. Then the multiple linear regression model takes the form

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon, \quad (1.3)$$

where X_j represents the j -th predictor and β_j quantifies the association between that variable and the response. We interpret β_j as the *average effect* on Y of a one-unit increase in X_j , holding all other predictors fixed.

1.2 Logistic Regression

Logistic regression is used to model the probability of a binary outcome Y based on predictor variables X . Unlike linear regression, which can yield probabilities outside the range $[0, 1]$, logistic regression uses the logistic function to constrain the predicted probabilities:

$$p(X) = \frac{e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \epsilon}}{1 + e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \epsilon}} = \frac{1}{1 + e^{-\beta_0 - \beta_1 X_1 - \cdots - \beta_p X_p - \epsilon}} \quad \text{where } p(X) = \Pr(Y = 1|X). \quad (1.4)$$

The function (1.4) produces an S-shaped curve, ensuring the probabilities are always between 0 and 1. The odds are defined as:

$$\text{Odds} = \frac{p(X)}{1 - p(X)}, \quad (1.5)$$

and taking the logarithm of the odds leads to the log-odds (logit):

$$\log \left(\frac{p(X)}{1 - p(X)} \right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon. \quad (1.6)$$

Here, β_1 represents the change in log-odds for a one-unit increase in X , and e^{β_1} is the factor by which the odds change. The logistic regression model is estimated using maximum likelihood estimation to ensure valid probability predictions.

Ultimately, logistic regression provides a flexible approach for modeling binary outcomes, ensuring valid probabilities through the logistic function and offering interpretability via odds and log-odds transformations.

2 Methodology

Linear Regression

To fit a multiple linear regression model using least squares, we use the `lm()` function in R. The model is specified as:

$$y \sim x_1 + x_2 + \cdots + x_n$$

where x_1, x_2, \dots, x_n are the predictors. The `summary()` function is then used to output the regression coefficients for all the predictors.

We can access individual components of the summary object by name. For instance, the R^2 , and `summary(lm.fit)$sigma` are given by `summary(lm.fit)$r.sq` and provides the residual standard error (RSE). Additionally, the `vif()` function from the `car` package is used to compute the variance inflation factors, which help to assess multicollinearity.

Logistic Regression

To fit a logistic regression model, we use the `glm()` function in R with the binomial family, which specifies that the response variable is binary. The model is expressed as:

$$\text{Response} \sim \text{Predictor}_1 + \text{Predictor}_2 + \cdots + \text{Predictor}_p$$

This model uses maximum likelihood estimation to estimate the coefficients for the predictors. After fitting the model, the `summary()` function provides the coefficients and their associated p-values, which indicate the significance of each predictor.

3 Model Evaluation and Interpretation

Linear Regression

In multiple linear regression, the primary objectives are to assess the relationship between the response variable Y and the predictors X_1, X_2, \dots, X_p , identify important predictors, evaluate model fit, and make predictions with quantifiable uncertainty.

- **Relationship Between Response and Predictors:** We test the null hypothesis $H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$ to determine if any predictors are related to Y . This is done using the F-statistic:

$$F = \frac{(\text{RegSS})/p}{\text{RSS}/(n - p - 1)} \quad (1.7)$$

A significant F-statistic (small p-value) should indicate that at least one predictor is useful.

- **Identifying Important Variables:** Variable selection methods, such as forward, backward, and mixed selection, are used to identify significant predictors, balancing model simplicity and predictive power.
- **Model Fit:** The model fit is evaluated using R^2 , which shows the proportion of variance explained, and the Residual Standard Error (RSE), which estimates model accuracy. The Adjusted R^2 accounts for the number of predictors.
- **Predictions and Uncertainty:** Predictions are made based on the fitted model, with uncertainty quantified by confidence intervals for coefficients and predictions. Model bias and irreducible error are considered as sources of prediction uncertainty.

Logistic Regression

- **Relationship Between Response and Predictors:** In logistic regression, the relationship between the response variable Y and the predictor variables is modeled using the logistic function. The significance of the predictors is evaluated using the Wald z-statistic, which tests the null hypothesis

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$$

If the p-values are small, it indicates that the corresponding predictor variables are significantly related to the outcome.

- **Identifying Important Variables:** In logistic regression, the significance of each predictor is assessed using the p-value associated with the coefficient. A low p-value (typically < 0.05) suggests that the predictor has a significant relationship with the response variable. Using variable selection techniques, such as backward elimination, could help identify the most important predictors.
- **Model Fit:** The model fit is evaluated using the Deviance and AIC (Akaike Information Criterion). The Deviance compares the fitted model with a baseline model (null model). Lower deviance values suggest a better model fit. The AIC helps in comparing different models, where a lower AIC value indicates a better fit while considering model complexity. Additionally, the Nagelkerke R^2 is a pseudo- R^2 measure that provides an indication of how well the model explains the variability of the response variable.
- **Checking Residuals:** Residuals in logistic regression, such as deviance residuals, are used to assess the goodness of fit and identify potential outliers or influential observations. Plots of residuals versus fitted values can be examined to check for patterns that may indicate model inadequacy.
- **Predictions and Uncertainty:** Predictions from the logistic regression model are made using the `predict()` function, which provides the probabilities of the event. These probabilities can be converted into class labels by setting a threshold. The uncertainty in predictions is measured through the standard errors of the coefficients. A model with higher accuracy will result in fewer misclassifications and a higher correct prediction rate.
- **Evaluating Predictions:** The performance of the logistic regression model can be assessed using a confusion matrix, which summarizes the counts of true positives, true negatives, false positives, and false negatives. Common metrics derived from the confusion matrix, such as accuracy, sensitivity, specificity, and precision, help evaluate the model's performance.

4 Data Pre-Processing

The dataset used in this analysis is the **tips dataset** from the **reshape2** package in R. This dataset provides information on restaurant tips, including variables such as the total bill, table size, day and time of day, and customer demographics (e.g., sex, smoking status). A log transformation was applied to the **tip** and **total.bill** variables to address potential skewness as seen in 1. The transformed variables are now referred to as **log.tip** and **log.total.bill**.

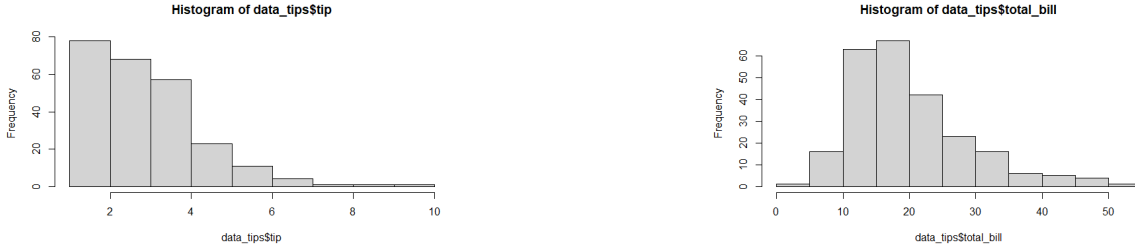


Figure 1: Histograms of Log Tip Amount and Log Total Bill

After the log transformation, the outliers were removed. This made the model have less explanatory power but in general both BIC and AIC were lower. The categorical variables (**size**, **sex**, **day**, **smoker** and **time**) were already encoded as factors, and the numerical variables (**log_tip**, **log_total_bill**) were in the correct format for analysis.

For the Logistic Regression part the response variable (**log_tip**) is going to be transformed, separating values greater than or equal to its median ('1') from the lower values ('0').

5 Exploratory Data Analysis

Statistical, Numerical, and Graphical Description of the Data

This section provides a statistical summary of the key variables in the dataset. The focus is on the log-transformed tip (**log_tip**), log-transformed total bill (**log_total_bill**), and table size (**size**).

Table 1: Summary Statistics of Key Variables

Variable	Mean	Standard Deviation	Range (Min–Max)
Log Tip (log_tip)	1.00	0.44	0.00–2.30
Log Total Bill (log_total_bill)	2.89	0.44	1.12–3.93
Size (size)	2.57	0.95	1.00–6.00

The categorical variables are summarized below:

Variable	Category	Count
Sex	Female	87
	Male	157
Smoker	Non-Smoker	151
	Smoker	93
Day	Sunday	19
	Thursday	87
	Tuesday	76
	Friday	62
Time	Dinner	176
	Lunch	68

Table 2: Summary of categorical variables.

The data reveals that log total bill values (mean = 2.89) are generally higher than log tip values (mean = 1.00), with significant variability (SD = 0.44 for both). The wide ranges suggest diverse tipping and bill behaviors. Table size shows moderate variation, while differences in categorical variables like sex and smoker status likely contribute to tipping differences.

Now we will see some graphical visualizations to explore the distribution of key variables and relationships in the dataset.



Figure 2: Histograms of Log Tip Amount and Log Total Bill

Figure 2 shows the histograms for the log-transformed tips (`log_tip`) with a vertical line representing the median and the log-transformed total bill (`log_total_bill`). Both histograms show a more balanced, approximately normal distribution centered around their respective means, reflecting fewer outliers.

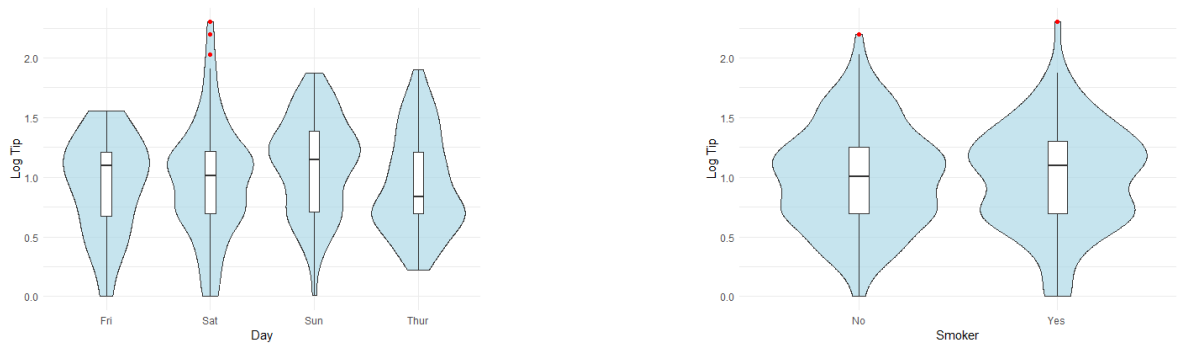


Figure 3: Violin Plots of Log Tip Amount by Day and Smoker

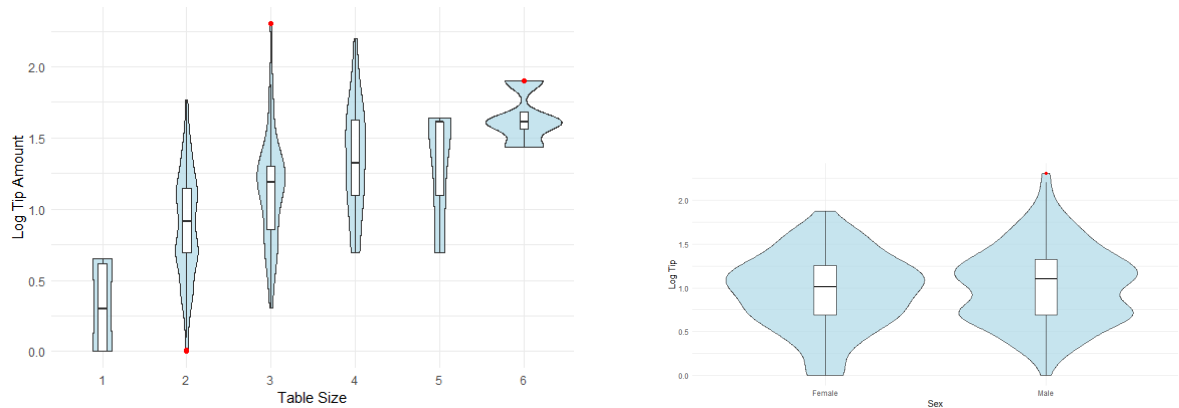
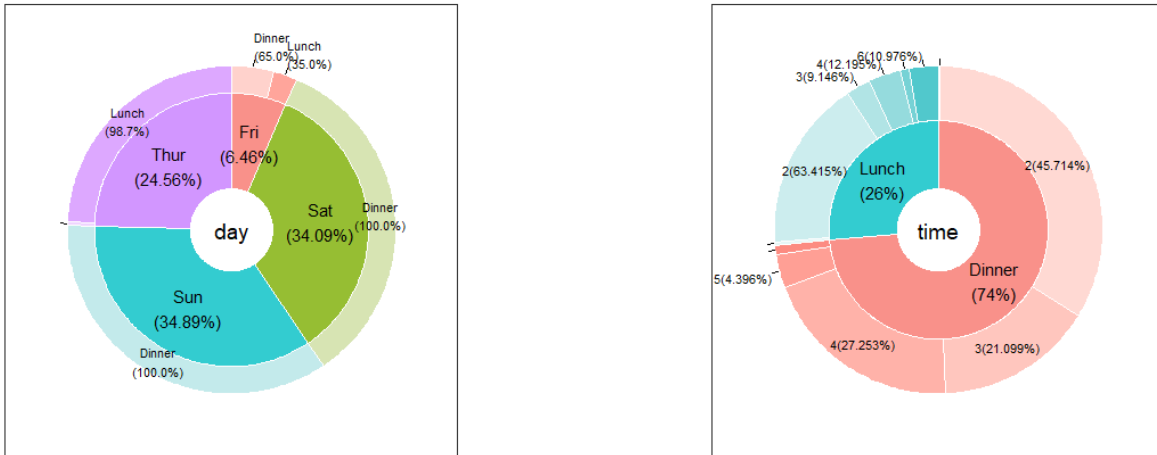
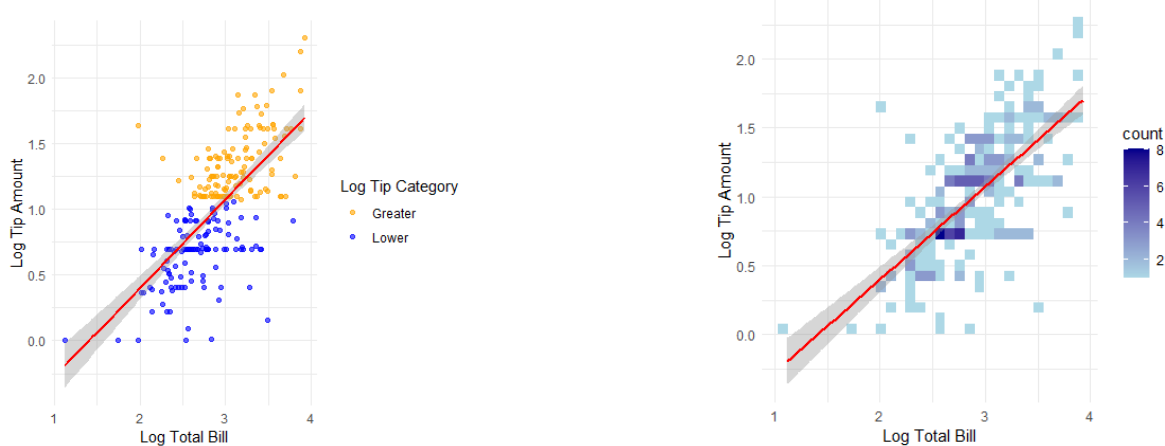


Figure 4: Violin Plots of Log Tip Amount by Size and Sex

Violin plots were used to explore the distribution of tips across different categorical variables. Figure 3 and figure 4 shows the distributions of `log_tip` by `day`, `smoker`, and table size (`size`). Larger table sizes generally show bigger in log tips. Both smoking and sex status have minimal impact on tipping. Across days, tipping is most variable on Saturdays, with more high-value outliers, while Thursdays have the most consistent behavior. The log transformation effectively normalizes tip amounts, highlighting trends across these features.



The first plot in Figure 5 illustrates the distribution of dining occurrences by day, highlighting that Saturday and Sunday dominate with approximately 35% each, followed by Thursday at 24% and Friday at a modest 6%. Lunch is prominent only on Thursday, while Friday exhibits the most balanced distribution between lunch and dinner. The second plot emphasizes the time distribution, revealing that dinner accounts for 74% of the total entries, compared to lunch at 26%. Additionally, the plots provide insight into group sizes at dinner, showing that groups of size 4, likely families, are more prevalent during this time. Overall, both plots suggest that weekends are the busiest days, with dinner being the preferred time for dining. Notably, no lunch entries are recorded on Saturdays and Sundays, possibly indicating the venue is closed during those times.



The first plot in Figure 6 visualizes the log tip amount categorized into lower and greater than the median, with the points colored accordingly and the regression line illustrating the relationship for both categories. The second plot shows a heatmap of the log-transformed total bill and tip amount, with a linear regression line (in red) highlighting the trend. Both plots imply that there is a correlation between this two variables.

6 Model Selection Process

The model selection process was carried out iteratively to identify the most suitable predictors for explaining the variation in the response variable (`log.tip`). Several models were evaluated using statistical metrics, as discussed in Section 3. The following is the step-by-step process:

Linear Model

Model	Formula	R^2	Adj. R^2	F-statistic	AIC	BIC
Initial w/o outliers	$\log(\text{tip}) \sim \log(\text{total.bill}) + \text{sex} + \text{day} + \text{smoker} + \text{size} + \text{time}$	0.448	0.428	23.48	140.5	175.3
Model 0	$\log(\text{tip}) \sim \log(\text{total.bill}) + \text{sex} + \text{day} + \text{smoker} + \text{size} + \text{time}$	0.478	0.460	26.9	147.9	182.9
Model 1	$\log(\text{tip}) \sim \log(\text{total.bill}) + \text{size} + \text{sex} + \text{time} + \log(\text{total.bill} * \text{size}) + \log(\text{total.bill} * \text{time})$	0.475	0.462	35.79	145.1	173.1
Model 2	$\log(\text{tip}) \sim \log(\text{total.bill}) + \text{size} + \text{time} + \log(\text{total.bill} * \text{size}) + \log(\text{total.bill} * \text{time})$	0.475	0.464	43.05	143.3	167.8
Model 3	$\log(\text{tip}) \sim \log(\text{total.bill}) + \text{size} + \text{time} + \log(\text{total.bill} * \text{size})$	0.472	0.465	71.46	140.8	162.7
Model 4	$\log(\text{tip}) \sim \log(\text{total.bill}) + \text{time} + \text{size}$	0.448	0.428	23.48	140.5	158.3
Final	$\log(\text{tip}) \sim \log(\text{total.bill}) + \text{size}$	0.4717	0.4674	107.6	138.8	152.8
Final w/o outliers	$\log(\text{tip}) \sim \log(\text{total.bill}) + \text{time} + \text{size} + \log(\text{total.bill} * \text{time})$	0.438	0.433	92.57	132.7	146.7

Table 3: Summary of Models with Key Metrics and AIC

We began with an initial model including all predictors without outliers. Then in Model 0 we kept outliers, evaluating significance and fit. Interaction terms were introduced in Model 1, capturing potential relationships between `log total bill`, `size`, and `time`. However, non-significant predictors such as `sex` and `smoker` were removed in Model 2 to improve parsimony. In Model 3, the interaction term `log total bill*time` was excluded due to lack of relevance, leaving only `log total bill`, `size`, and their interaction. Model 4 further simplified the predictors by excluding interaction terms, retaining only `log total bill`, `size`, and `time`. Finally, the **Final** model focused on the essential variables (`log total bill` and `size`) after reintroducing log-transformation, as the relationship between tip and total bill is typically multiplicative rather than additive.

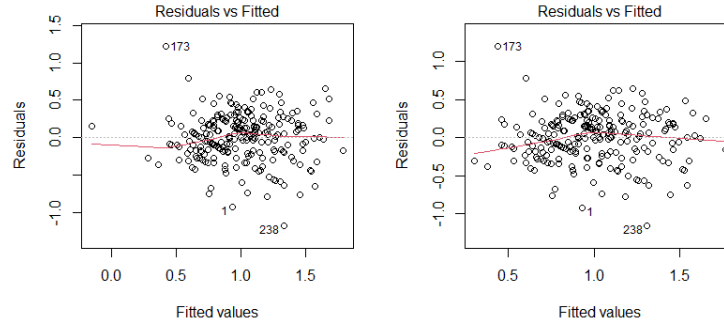


Figure 7: Residuals with vs without outliers

Throughout this process, model refinement aimed to balance interpretability and statistical rigor. In Figure 7 we see that the behavior of residuals is similar with and without outliers. That being, the final model without outliers, which included only `log total bill` and `size`, demonstrated strong explanatory power with the lowest AIC (132.7) and BIC (146.7) values between all models, which confirms it as the most effective representation of tip variations.

Logistic Regression Models

Model	Formula	R^2_{MF}	Resid Deviance	Nagelkerke	AIC	BIC
Initial w/o outliers	$\log(\text{tip.category}) \sim \log(\text{total.bill}) + \text{size} + \text{sex} + \text{smoker} + \text{time} + \text{day}$	0.300	233.96	0.453	251.96	283
Initial	$\log(\text{tip.category}) \sim \log(\text{total.bill}) + \text{size} + \text{sex} + \text{smoker} + \text{time} + \text{day}$	0.303	235.68	0.458	253.68	285
Model 1	$\log(\text{tip.category}) \sim \log(\text{total.bill}) + \text{size} + \text{sex} + \text{time} + \text{day}$	0.302	235.98	0.457	251.98	280
Model 2	$\log(\text{tip.category}) \sim \log(\text{total.bill}) + \text{size} + \text{time} + \text{day}$	0.301	236.38	0.455	250.38	275
Model 3	$\log(\text{tip.category}) \sim \log(\text{total.bill}) + \text{time} + \text{day}$	0.300	237.05	0.453	249.05	270
Model 4	$\log(\text{tip.category}) \sim \log(\text{total.bill}) + \text{time}$	0.285	242.00	0.435	247.45	258
Model 5	$\log(\text{tip.category}) \sim \log(\text{total.bill}) * \text{time}$	0.299	237.04	0.453	245.05	259
Final w/o outliers	$\log(\text{tip.category}) \sim \log(\text{total.bill}) * \text{day}$	0.320	226.86	0.479	242.86	271
Final	$\log(\text{tip.category}) \sim \log(\text{total.bill}) * \text{day}$	0.328	227.14	0.488	243.14	271

Table 4: Summary of Logistic Regression Models with Key Metrics, AIC, and BIC

We began with an initial logistic regression model that included all predictors: `log_total_bill`, `size`, `sex`, `smoker`, `time`, and `day`. The `smoker` variable was first removed as it was not statistically significant. Next, the `sex` variable was excluded for similar reasons, leaving `log_total_bill`, `size`, `time`, and `day`. The model was further simplified by removing `size`, leaving only `log_total_bill`, `time`, and `day`. Subsequently, the model was reduced to include only `log_total_bill` and `time`. An interaction term between `log_total_bill` and `time` was introduced, but this did not yield significant improvements. We also explored using a polynomial term for `log_total_bill`, which showed similar performance to the previous models. Finally, the model focused solely on the interaction between `log_total_bill` and `day`, which resulted in the best overall performance, with the highest R^2 and Nagelkerke, lowest residual deviance, and almost the lowest AIC.

Throughout the model refinement process, each step aimed to reduce complexity while maintaining significant predictors. The final model demonstrated superior performance, highlighting the importance of interactions between `log_total_bill` and `day`. Despite these improvements, the model still is not the best, likely due to the simplification of the response variable into a binary category, which may not fully capture the variability in the data. Considering that the null deviance is 338.8, the residual deviance of 226.86 in the final model is satisfactory, as it follows the rule of thumb because $226.86 \approx 232 = K - (p+1)$, being $K = 240$ and $p = 7$. Also, as we can see in Figure 8 the residuals behave like we would expect they would. In Table 5 we can confirm everything we said about the model quality until now. With a 77.6% accuracy it is good but not great.

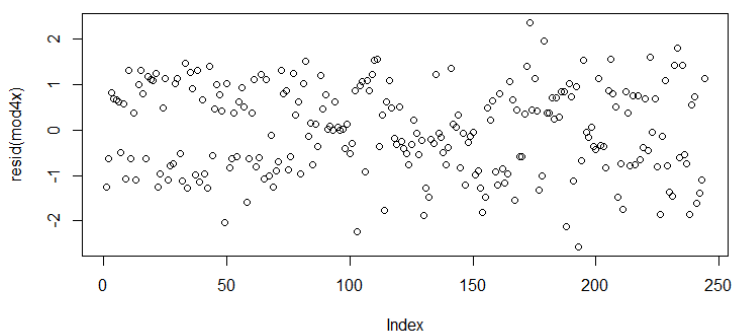


Figure 8: Residuals with vs without outliers

Prediction	0	1
0	87	20
1	34	100

Table 5: Confusion Matrix

7 Analysis of Variables

Linear Regression: Continuous Variable (`log_total_bill`) and Categorical Variable (`size`)

Raw and Adjusted Effects of `log_total_bill` and `size`

The raw effect of `log_total_bill` on the response variable is significant, with an estimated coefficient of 0.675 with a very high significance level, indicating that as `log_total_bill` increases by one unit, the response variable increases by approximately 67.5%. While the adjusted effect of `log_total_bill`, controlling for other variables, is 0.602 and remains significant.

As we can see in Table 6 the raw effects of `size` categories on the response variable are significant, with increasing magnitudes across categories, indicating meaningful associations. However, after ad-

Category	Raw Estimate	Raw p-value	Adjusted Estimate	Adjusted p-value
2	0.5640	0.00386	0.05346	0.754
3	0.8271	5.19e-05	0.11170	0.542
4	1.0315	5.88e-07	0.17584	0.353
5	1.0154	9.70e-05	0.11860	0.611
6	1.3249	1.70e-06	0.34020	0.170

Table 6: Raw and Adjusted Effects for size

justing for other variables, none of the adjusted effects are statistically significant, suggesting that the raw associations may be influenced by confounding factors. While Category 6 shows the largest raw effect (1.3249), the adjusted effect is much smaller and non-significant ($p = 0.170$). These findings imply that size categories are not independent predictors of the response variable after accounting for other factors in the model. Additionally, this suggests that a linear model may not be the best approach for this dataset.

Effect of Changing Categories in Size

The effect of changing from the third category of size to the second category is estimated to be -0.058 , indicating a small decrease in the response variable. However, the 95% confidence interval for this effect is $[-0.548, 0.431]$, and the 90% confidence interval is $[-0.469, 0.353]$. Since both confidence intervals include zero, we conclude that the change from the third to the second size category is not statistically significant.

Investigation of Interaction Between Log Total Bill and Size

The interaction terms between `log_total_bill` and `size` were found to be statistically insignificant (all $p > 0.05$), suggesting that the relationship between `log_total_bill` and `log_tip` does not substantially vary across different size categories. This is further confirmed by an ANOVA comparison, which shows no significant improvement in model fit when including interaction terms ($p = 0.551$). Additionally, the adjusted R^2 (0.4577) and the residual standard error (0.3212) indicate minimal differences between models with and without the interaction terms. Thus, the simpler additive model, without interactions, is sufficient to describe the data and provides a more straightforward interpretation without sacrificing explanatory power.

Logistic Regression: Continuous Variable (`log_total_bill`), Categorical Variable (`day`)

Raw and Adjusted Effects of Log Total Bill and Day of the Week

The raw effect of log total bill on the likelihood of receiving a tip is highly significant, with an estimated coefficient of 3.9430 ($p < 0.001$). This result suggests that for every one-unit increase in the log total bill, the odds of receiving a tip increase substantially. The adjusted effect of log total bill remains equally strong and significant, with a similar coefficient of 3.3376 ($p < 0.001$). This confirms a robust positive relationship between the total bill and the likelihood of tipping, even after accounting for the day of the week.

Category	Raw Estimate	Raw p-value	Adjusted Estimate	Adjusted p-value
Saturday	-0.5351	0.123	-1.4708	0.0598
Sunday	-0.3245	0.215	-0.5021	0.311
Thursday	-0.2504	0.315	-0.1320	0.511

Table 7: Raw and Adjusted Effects for Day of the Week

As shown in Table 7, the raw effects for day of the week are not significant. After adjusting for the other variables on the model, only the effect for Saturday approaches significance. This suggests that tipping behavior on Saturday may slightly differ from Friday, though the evidence is marginal. The effects of Sunday and Thursday remain non-significant in the adjusted model, indicating no substantial differences in tipping likelihood on these days compared to Friday.

Effect of Changing Categories in Day of the Week

The effect of changing from the third category of **day** (Sunday) to the second category (Saturday) is estimated to be 0.0864, indicating a small increase in the likelihood of receiving a tip. However, the 95% confidence interval for this effect is $[-2.4247, 2.5974]$, and the 90% confidence interval is $[-2.0211, 2.1939]$. Since both confidence intervals include zero, we conclude that the change from Sunday to Saturday is not statistically significant.

Investigation of Interaction Between Log Total Bill and Day of the Week

The interaction terms between log total bill and day of the week were found to be statistically significant for the model comparison ($p = 0.01063$), suggesting that the relationship between log total bill and the likelihood of receiving a tip does indeed vary across different days of the week.

The coefficients for Saturday, Sunday, and Thursday are -10.190, -9.642, and -5.468, respectively, with p-values of 0.1304, 0.1534, and 0.4349. These results indicate no significant effects for these days compared to Friday.

However, the ANOVA comparison between the models with and without the interaction terms yields a p-value of 0.01063, which is less than 0.05. This indicates that the inclusion of interaction terms significantly improves the model fit, suggesting that the effect of log total bill on the likelihood of receiving a tip does vary depending on the day of the week, justifying the inclusion of the interaction terms in the model. Additionally, the Akaike Information Criterion (AIC) for the model with interaction is 243.14, compared to 248.35 for the simpler additive model.

The inclusion of interaction terms between log total bill and day of the week improves model fit, as indicated by the ANOVA ($p = 0.01063$) and the lower residual deviance and AIC. Although the coefficients for individual days were not significant, the interaction model offers a better fit overall and is preferred over the simpler model.

8 Conclusion

The linear regression analysis demonstrated that the log-transformed total bill and table size are the most important predictors of tipping behavior. The model refinement process revealed that including interaction terms did not significantly improve the model's performance. The final model, which used only log-transformed total bill and size as predictors without the outliers, had the lowest AIC and provided a straightforward explanation of the data. However, adjusted effects for table size categories were not significant, suggesting confounding factors influence the raw associations. Interaction terms between log total bill and size were also insignificant, confirming that a simpler additive model was sufficient. The linear model highlights a strong positive association between log total bill and tip amounts but suggests limited predictive power for other variables like sex, smoker status, and day of the week.

The logistic regression analysis explored the relationship between some features and above average tipping likelihood, categorized into high or low tips. The final logistic model, which incorporated interaction terms between log total bill and day of the week, showed the best fit with the lowest AIC and highest Nagelkerke and McFadden R^2 . While individual effects for day categories were not significant, the interaction terms significantly improved model fit, indicating that tipping likelihood varies across days depending on the total bill. Even though it was a satisfactory model, the binary classification approach may have oversimplified the variability in tipping behavior, potentially limiting model performance. Nonetheless, the results underscore the robust positive influence of total bill on tipping likelihood, with nuanced variations by day.

To conclude, the linear model effectively captures the continuous relationship between tipping amounts and predictors, while the logistic model provides insights into tipping likelihood and its interactions with categorical variables. The linear model is preferred for detailed predictions of tip amounts, whereas the logistic model highlights trends in tipping behavior. Both models underscore the importance of log-transformed total bill as the dominant predictor but suggest limited utility for other factors like table size, sex, and smoker status.

References

- [1] An Introduction to Statistical Learning. Available: <https://www.statlearning.com/>
- [2] Classification Metrics Walkthrough. [Online]. Available: <https://www.kdnuggets.com/2022/10/classification-metrics-walkthrough-logistic-regression-accuracy-precision-recall-roc.html>
- [3] Linear and Logistic Regression Slides from Moodle.