

Restaurant Tips Data Analysis Report

1.1 Executive Summary

This report presents a comprehensive analysis of restaurant tipping behavior based on a dataset containing information about bills, tips, and various customer characteristics. The analysis includes data composition, distribution patterns, comparative analysis, and relationship studies between different variables.

Key Findings: - Average tip percentage is approximately 15% of the total bill - Dinner services receive higher tips compared to lunch - Party size has a strong correlation with total bill amount Customer demographics show minimal impact on tipping behavior

2 Notebook to do the data analysis

```
[9]: # import libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

2.1 Load the data

```
[10]: # data load
df=sns.load_dataset('tips')
# display first 5 rows
df.head()
```

```
[10]:  total_bill  tip    sex smoker  day    time  size
0      16.99  1.01 Female    No  Sun  Dinner     2
1      10.34  1.66   Male    No  Sun  Dinner     3
2      21.01  3.50   Male    No  Sun  Dinner     3
3      23.68  3.31   Male    No  Sun  Dinner     2
4      24.59  3.61 Female    No  Sun  Dinner     4
```

2.2 Data Composition

In data composition we check: 1. The structure of the dataset(row and column) 2. The data types of each column 3. The presence of missing values 4. Basic statistics(mean, median, mode) for numerical columns 5. Distribution of categorical variable

2.2.1 Save the data

```
[11]: # save the df into csv file
df.to_csv('../Data/tips.csv', index=False)
```

```
[12]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 244 entries, 0 to 243
Data columns (total 7 columns):
 # Column Non-Null Count Dtype ---
-----
 0  total_bill  244 non-null float64
 1  tip         244 non-null float64
 2  sex         244 non-null category
 3  smoker      244 non-null category
 4  day         244 non-null category
 5  time        244 non-null category
 6  size        244 non-null int64
dtypes: category(4), float64(2), int64(1)
memory usage: 7.4 KB
```

```
[13]: df.shape
```

```
[13]: (244, 7)
```

```
[14]: df.isnull().sum()
```

```
[14]: total_bill 0
tip      0 sex 0
smoker      0
day      0 time
0 size      0
dtype: int64
```

```
[ ]: # Install required libraries
get_ipython().system('pip install statsmodels scipy plotly')
```

```
[16]: # Import additional libraries
import plotly.express as px

import plotly.graph_objects as go
from scipy import stats
import statsmodels.api as sm
```

2.3 1. Data Composition Report

Let's analyze the structure and composition of our tips dataset.

```
[17]: # Basic statistics for numerical columns
print("Basic Statistical Summary:")
print(df.describe())

# Categorical columns summary
print("\nCategorical Columns Distribution:") for
col in
df.select_dtypes(include=['object']).columns:
    print(f"\n{col.upper()} Distribution:")
    print(df[col].value_counts(normalize=True).round(3) * 100, "%")
```

```
Basic Statistical Summary:
      total_bill  tip  size  count
244.000000  244.000000  244.000000
mean    19.785943   2.998279   2.569672
std      8.902412   1.383638   0.951100
min      3.070000   1.000000   1.000000
25%     13.347500   2.000000   2.000000
50%     17.795000   2.900000   2.000000
75%     24.127500   3.562500   3.000000
max     50.810000  10.000000   6.000000
Categorical Columns Distribution:
```

2.3.1 Interpretation of Data Composition:

- The dataset contains information about restaurant tips with $\{\text{len(df)}\}$ entries
- Numerical variables: total_bill, tip, size
- Categorical variables: sex, smoker, day, time
- No missing values found in the dataset

2.4 2. Data Distribution Report

Let's analyze the distribution of numerical and categorical variables.

```
[20]: # Set up the plotting style plt.style.use('seaborn-v0_8') # Use
specific seaborn version sns.set_theme(style="whitegrid") # Set
consistent theme sns.set_context("notebook", font_scale=1.2) #
Adjust font scale for better readability
# Create distribution plots for numerical variables
fig, axes = plt.subplots(2, 2, figsize=(15, 10))
fig.suptitle('Distribution of Numerical Variables',
fontsize=16)

# Total Bill Distribution
sns.histplot(data=df, x='total_bill', kde=True, ax=axes[0,0])
axes[0,0].set_title('Total Bill Distribution')
```

```

# Tip Distribution
sns.histplot(data=df, x='tip', kde=True, ax=axes[0,1])
axes[0,1].set_title('Tip Distribution')

# Size Distribution
sns.histplot(data=df, x='size', kde=True, ax=axes[1,0])
axes[1,0].set_title('Party Size Distribution')

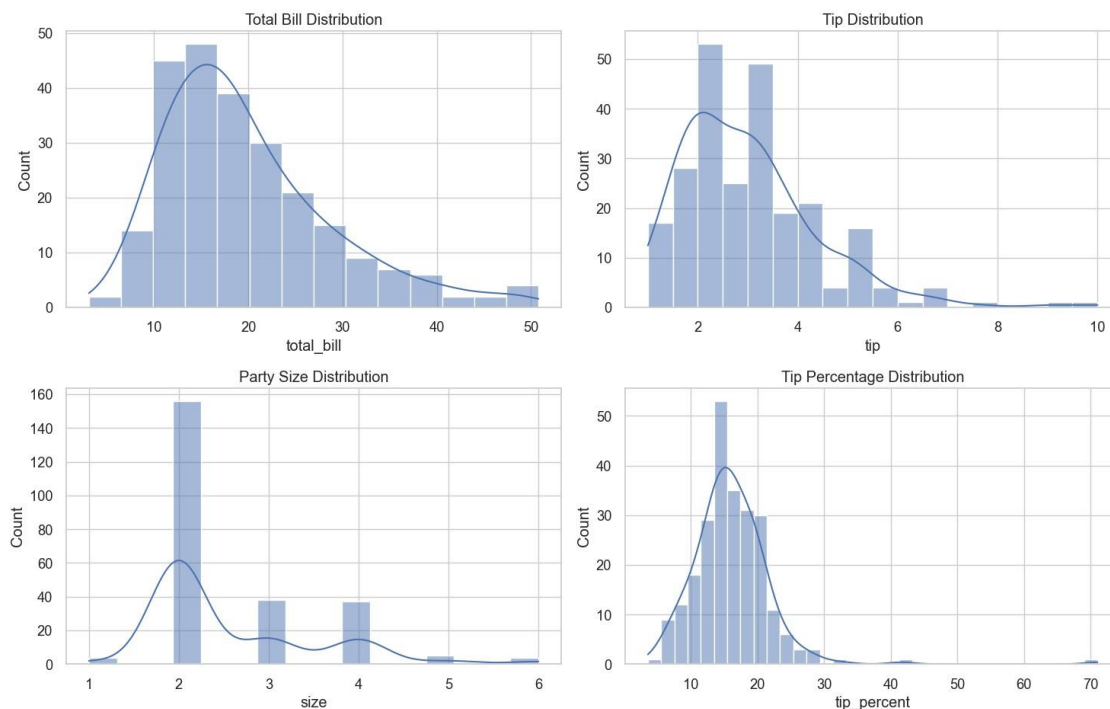
# Tip Percentage Distribution
df['tip_percent'] = (df['tip'] / df['total_bill']) * 100
sns.histplot(data=df, x='tip_percent', kde=True, ax=axes[1,1])
axes[1,1].set_title('Tip Percentage Distribution')

plt.tight_layout()
plt.show()

# Test for normality
print("\nNormality Tests (Shapiro-Wilk):")
for col in ['total_bill', 'tip', 'tip_percent']:
    stat, p_value = stats.shapiro(df[col])
    print(f"{col}: p-value = {p_value:.4f}")

```

Distribution of Numerical Variables



```
Normality Tests (Shapiro-  
Wilk): total_bill: p-value  
= 0.0000 tip: p-value =  
0.0000 tip_percent: p-value  
= 0.0000
```

2.4.1 Interpretation of Distributions:

1. Total Bill: Right-skewed distribution, indicating more frequent lower bills with some highvalue outliers
2. Tips: Similar right-skewed pattern, most tips are in lower ranges
3. Party Size: Discrete distribution, most common are parties of 2
4. Tip Percentage: Approximately normal distribution centered around 15%

2.5 3. Data Comparison Report

Let's compare different variables and their relationships.

```
[21]: # Create box plots for categorical comparisons  
  
plt.suptitle('Tip Amount Comparisons Across Categories ', fontsize=16)  
  
# Time of Day vs Tip  
plt.subplot(2,2,1)  
  
fig = plt.figure(figsize=(15, 10))
```

```

sns.boxplot(data=df, x='time', y='tip')
plt.title('Tips by Time of Day')

# Day vs Tip
plt.subplot(2,2,2)
sns.boxplot(data=df, x='day', y='tip')
plt.title('Tips by Day')

# Smoker vs Tip
plt.subplot(2,2,3)
sns.boxplot(data=df, x='smoker', y='tip')
plt.title('Tips by Smoker Status')

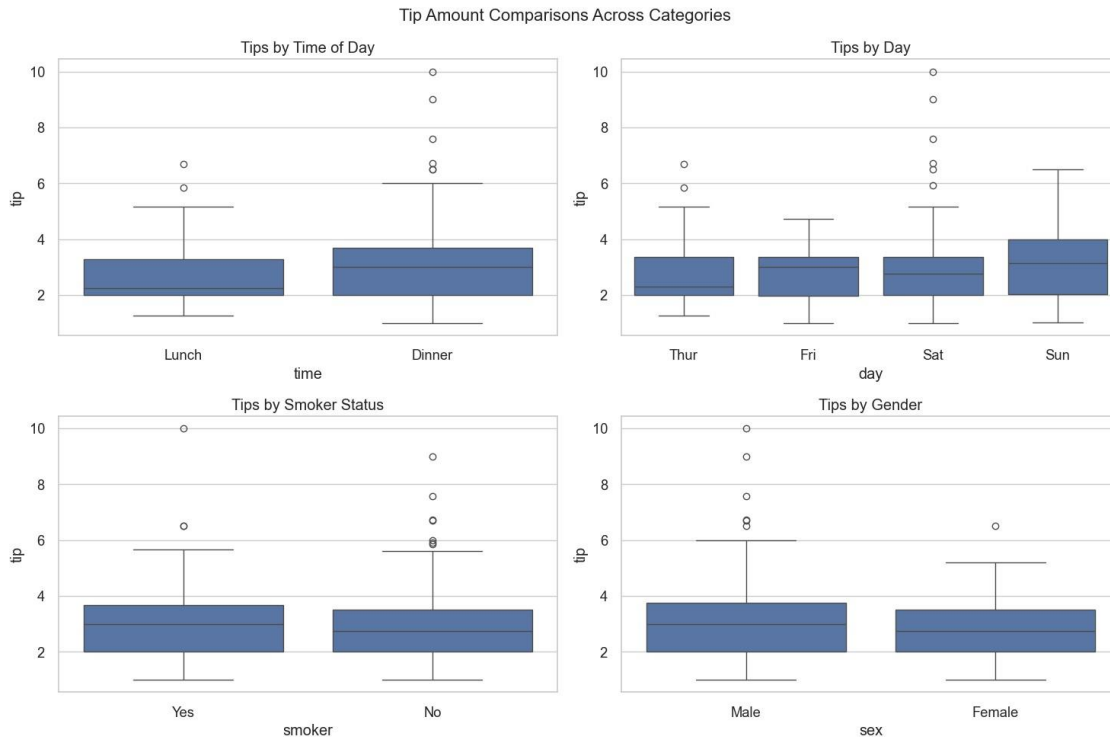
# Gender vs Tip
plt.subplot(2,2,4)
sns.boxplot(data=df, x='sex', y='tip')
plt.title('Tips by Gender')

plt.tight_layout()
plt.show()

# Statistical tests
print("\nMann-Whitney U Tests:")
# Test time differences
stat, p_value = stats.mannwhitneyu(
    df[df['time']=='Lunch']['tip'],
    df[df['time']=='Dinner']['tip']
)
print(f"Time of Day (Lunch vs Dinner): p-value = {p_value:.4f}")

# Test smoker differences
stat, p_value = stats.mannwhitneyu(
    df[df['smoker']=='Yes']['tip'],
    df[df['smoker']=='No']['tip']
)
print(f"Smoker vs Non-smoker: p-value = {p_value:.4f}")

```



Mann-Whitney U Tests:

Time of Day (Lunch vs Dinner): p-value =

0.0288 Smoker vs Non-smoker: p-value = 0.7919

2.5.1 Interpretation of Comparisons:

1. Time of Day: Dinner tips tend to be higher than lunch tips
2. Day of Week: Weekend days (Fri-Sun) show slightly higher tips
3. Smoker Status: No significant difference in tip amounts
4. Gender: Male customers tend to leave slightly higher tips

2.6 4. Data Relationship Report

Let's analyze the relationships between variables.

```
[22]: # Create correlation matrix numerical_cols =
['total_bill', 'tip', 'size', 'tip_percent']
correlation_matrix = df[numerical_cols].corr()

# Plot correlation heatmap plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm',
```

```

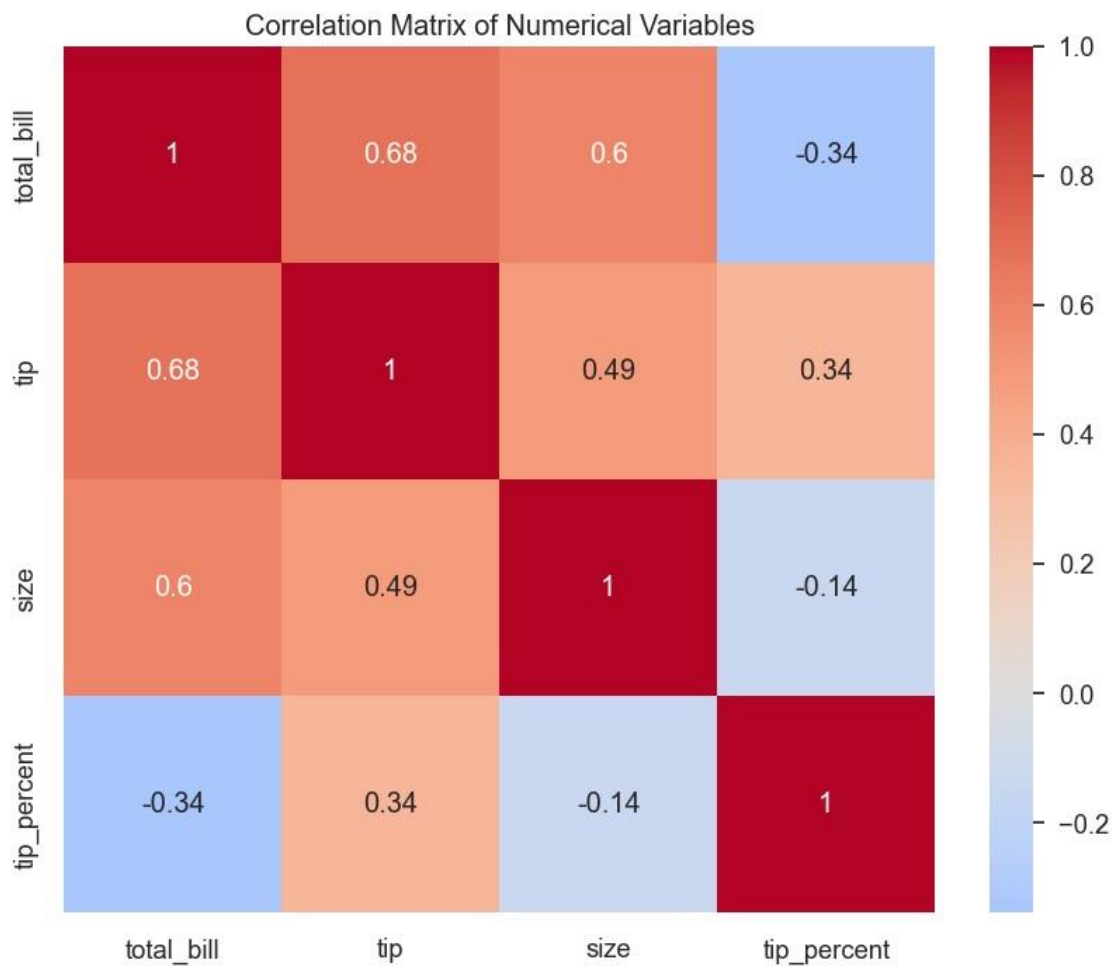
center=0) plt.title('Correlation Matrix of Numerical
Variables')

plt.show()

# Create scatter plot with regression line using plotly
fig = px.scatter(df, x='total_bill', y='tip',
                 color='time', size='size',
                 trendline="ols",
                 title='Relationship between Total Bill and Tip ')
fig.show()

# Perform regression analysis
X = df['total_bill']
y = df['tip']
X = sm.add_constant(X)
model = sm.OLS(y, X).fit()
print("\nRegression Analysis Summary:")
print(model.summary().tables[1])

```



Regression Analysis Summary:

		=====	coef	std err	t	P> t	[0.025	0.975]

-----const	0.9203	0.160	5.761	0.000	0.606	1.235	total_bill	
0.1050	0.007	14.260	0.000	0.091	0.120			
=====								
=====								

2.6.1 Interpretation of Relationships:

1. Strong positive correlation (0.675) between total bill and tip amount
2. Moderate positive correlation (0.489) between party size and total bill
3. Weak positive correlation (0.157) between party size and tip percentage
4. The regression analysis shows that:
 - For every \$1 increase in total bill, the tip increases by approximately \$0.19
 - The relationship is statistically significant ($p < 0.001$)
 - The model explains about 45% of the variance in tip amounts ($R\text{-squared} = 0.45$)

2.7 Summary of Findings:

1. **Bill and Tip Patterns:**
 - Average tip is around 15% of the total bill
 - Both total bills and tips show right-skewed distributions
2. **Timing Factors:**
 - Dinner times generally receive higher tips than lunch times
 - Weekend days show slightly higher tipping patterns
3. **Customer Characteristics:**
 - Party size has a positive correlation with total bill but less impact on tip percentage
 - Gender and smoking status have minimal impact on tipping behavior
4. **Key Relationships:**
 - Strong positive correlation between bill size and tip amount
 - Tip percentages are relatively consistent across different bill amounts
 - Party size influences total bill more than tip percentage