# Order Delivery Time Prediction: A Linear Regression Analysis

## 1. Project Objective

The primary goal of this project is to develop a robust regression model to accurately predict the delivery time for orders placed through Porter. By analyzing a wide range of features—including order specifics, restaurant details, and the availability of delivery partners—the model aims to uncover the key drivers of delivery duration. The insights from this model are intended to optimize logistical efficiency, improve customer satisfaction by setting accurate delivery expectations, and enhance the overall delivery service.

## 2. Data Understanding and Pipeline

The analysis was performed on the `porter_data_1.csv` dataset, which contains transactional order data. The project followed a structured data science pipeline:

1. Data Loading: Importing the raw dataset.
2. Data Preprocessing & Feature Engineering: Cleaning the data and creating new, more informative features.
3. Exploratory Data Analysis (EDA): Visually inspecting the data to uncover patterns and relationships.

4. Model Building: Developing the predictive model using selected features.
5. Model Evaluation & Inference: Assessing the model's performance and saving it for future use.

## 3. Data Preprocessing and Feature Engineering

A series of transformations were applied to prepare the dataset for modeling. These steps were crucial for handling inconsistencies and extracting maximum value from the raw data.

- Target Variable Creation: The core target variable, `time_taken` (in minutes), was engineered by calculating the difference between `actual_delivery_time` and `created_at`.
- Temporal Feature Extraction: To capture time-based patterns, the `created_at` timestamp was used to derive new features:
    - `order_placed_hour`: The hour of the day the order was placed.
    - `order_placed_day`: The day of the week.
    - `isWeekend`: A binary flag indicating if the order was placed on a weekend.
- Handling Missing Values: Rows with missing values in key operational fields like `total_onshift_dashers`, `total_busy_dashers`, and `store_primary_category` were removed to ensure data quality.
- Categorical Encoding: The `store_primary_category` feature was converted into numerical format using one-hot encoding, creating separate binary columns for each category.
- Numerical Scaling: All numerical features were standardized using `StandardScaler`. This process scales the data to have a mean of 0 and a standard deviation of 1, preventing features

with larger scales from unduly influencing the model.

- Feature Dropping: Redundant or unnecessary columns, including original timestamps and IDs like `market_id`, were dropped to simplify the model.

## 4. Exploratory Data Analysis (EDA)

EDA revealed several important characteristics of the dataset:

- Feature Distributions: The distributions of numerical features like `total_items` and `subtotal` were found to be right-skewed, indicating that most orders contain fewer items and have lower subtotals, with a smaller number of large orders.
- Time-Based Trends: Analysis showed that delivery times fluctuate significantly based on the hour of the day, with noticeable peaks during typical meal times.
- Correlations: A correlation matrix was used to examine relationships between variables, helping to identify potential multicollinearity before model building.

## 5. Model Development

A Linear Regression model was selected for its interpretability and efficiency.

- Feature Selection: To build a parsimonious and effective model, Recursive Feature Elimination (RFE) was employed. RFE iteratively removes the least important features, and based on this analysis, the top 9 features were selected. These included `total_items`, `subtotal`, and specific one-hot encoded hour features like `order_placed_hour_5` (5 AM) and `order_placed_hour_16` (4 PM).

- Model Training: The dataset was split into a training set (80%) and a testing set (20%). The Linear Regression model **was** trained on the standardized training data using the 9 features identified by RFE.

## 6. Results and Model Evaluation

The model's performance was assessed on the unseen test data.

- Performance Metric: The model achieved an R-squared ($R^2$) score of 0.1267. This means the model explains approximately 12.7% of the variance in delivery time. While this is a modest score, it indicates the selected features have some predictive power. The remaining variance is likely due to factors not included in the dataset, such as real-time traffic, weather, or specific order complexities.
- Residual Analysis:
  - Residuals vs. Predicted Plot: A plot of the model's residuals (errors) against the predicted values showed a random scatter of points around the zero line. This is a positive diagnostic result, suggesting that the assumption of linearity is appropriate and there is no significant heteroscedasticity (uneven error variance).
  - Histogram and Q-Q Plot: These plots confirmed that the residuals were approximately normally distributed, satisfying another key assumption of linear regression.
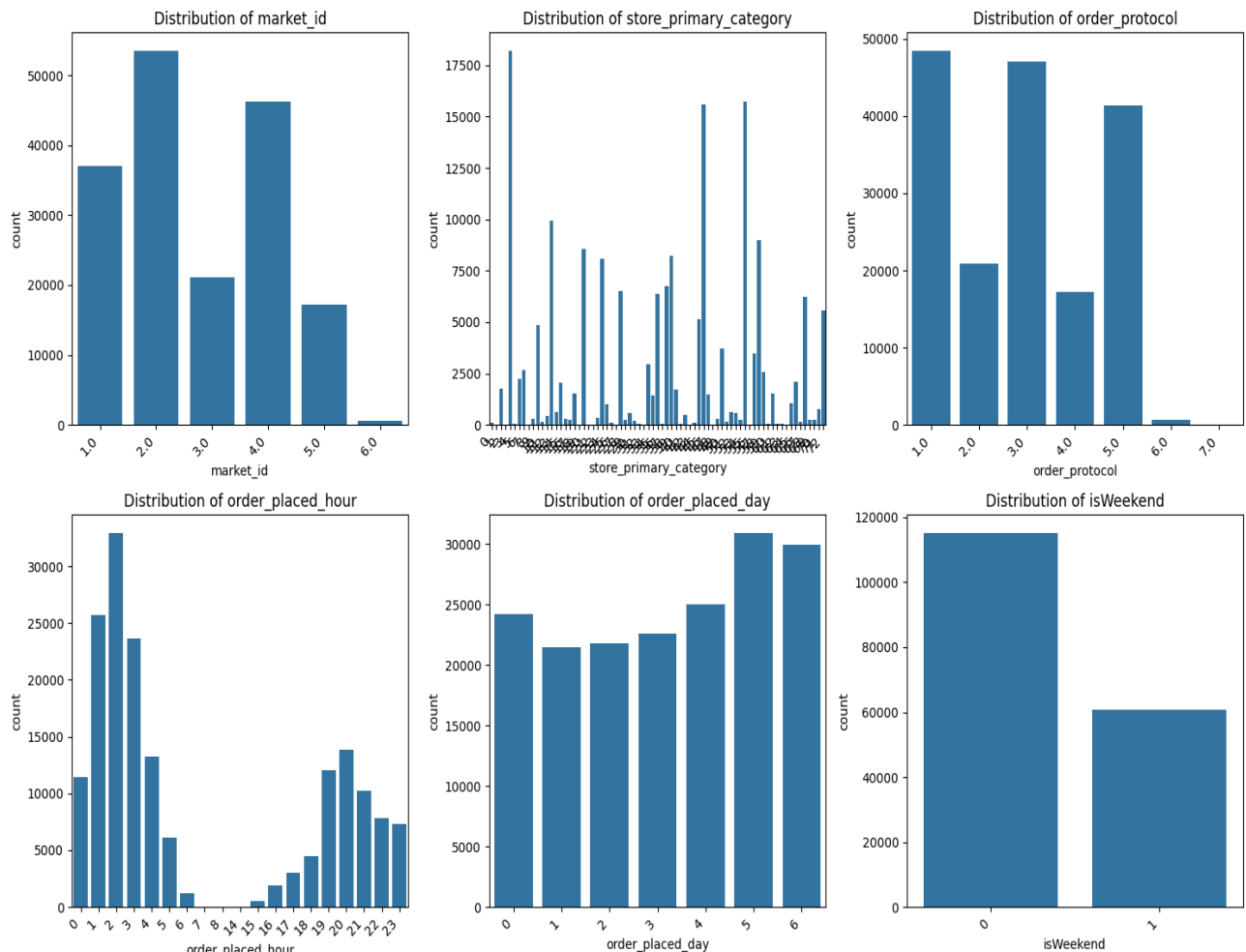
## 7. Insights and Business Outcomes

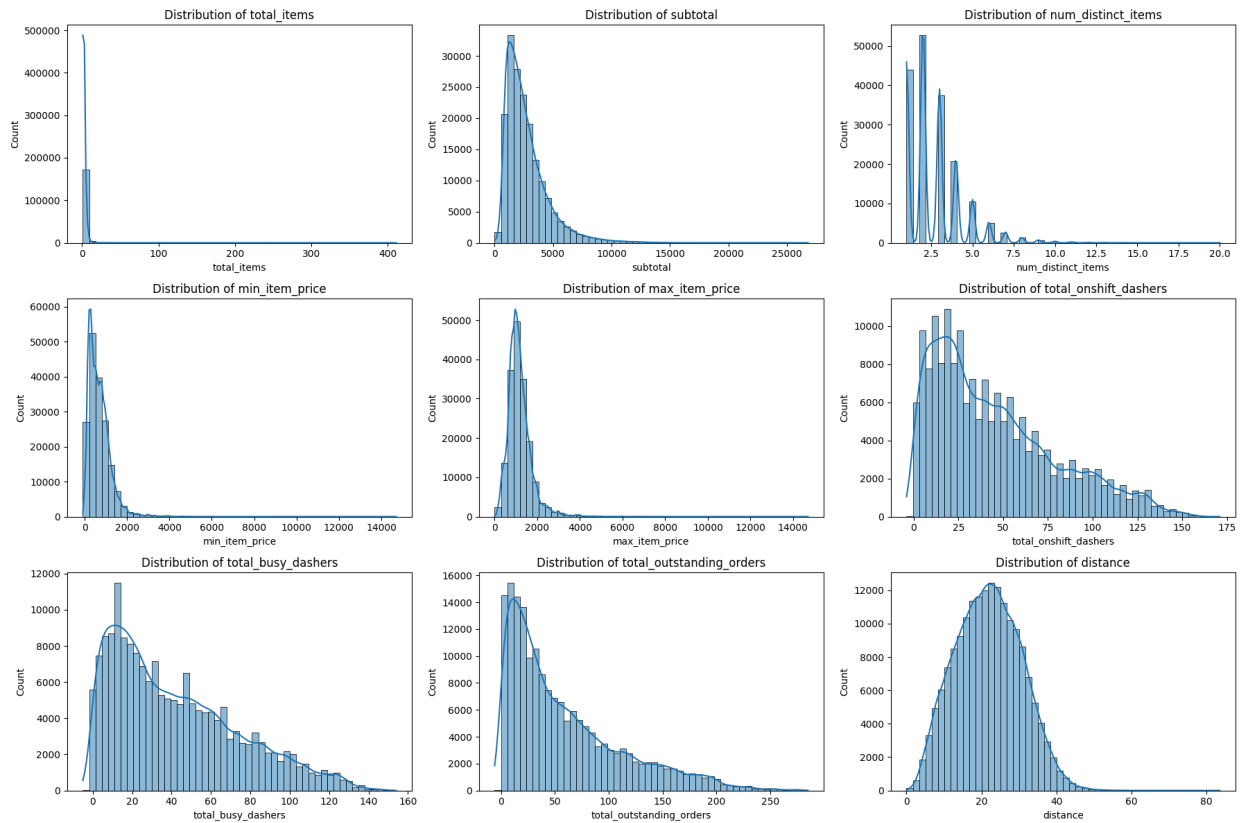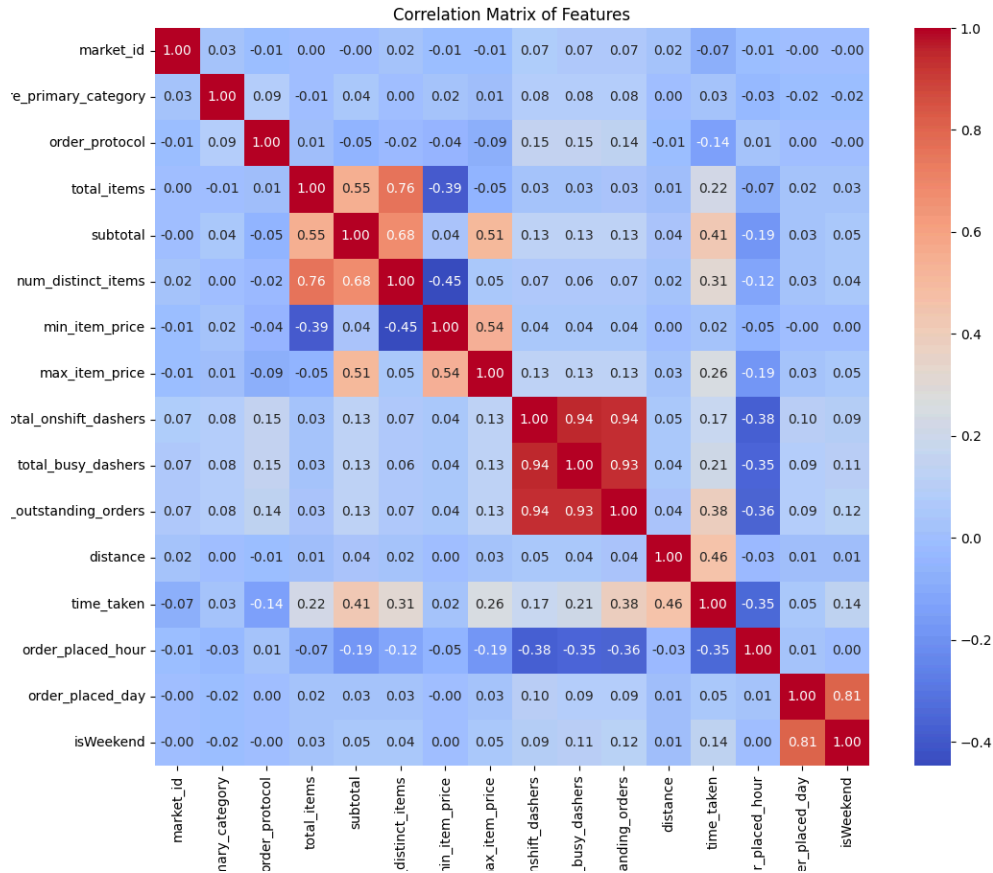The analysis and resulting model provide several actionable insights:

- Key Drivers of Delivery Time: The most influential predictors
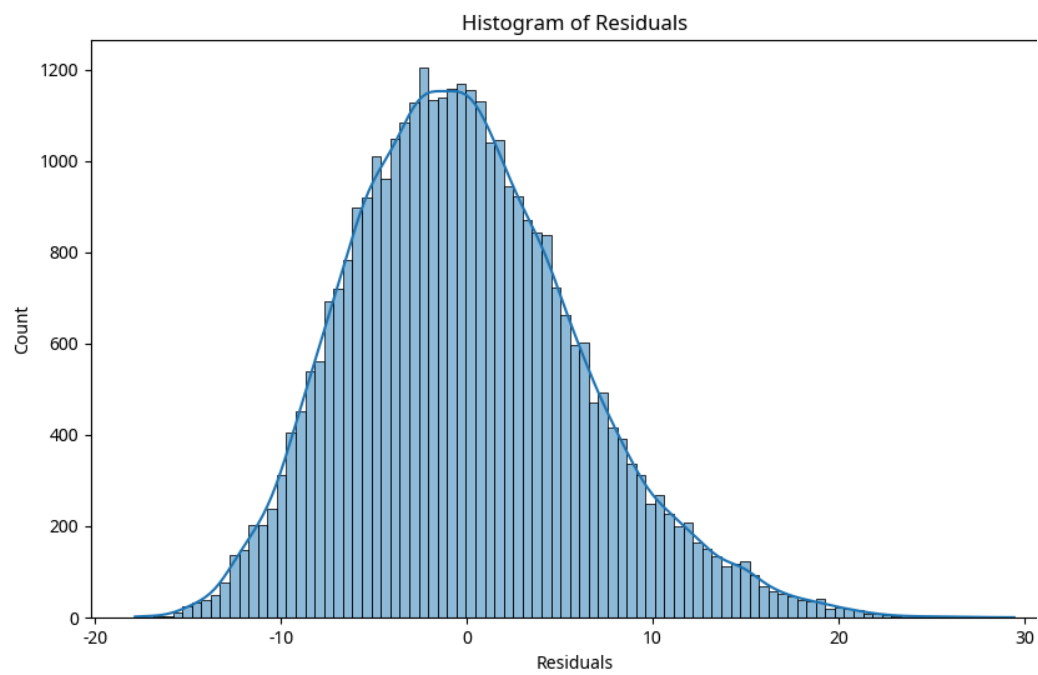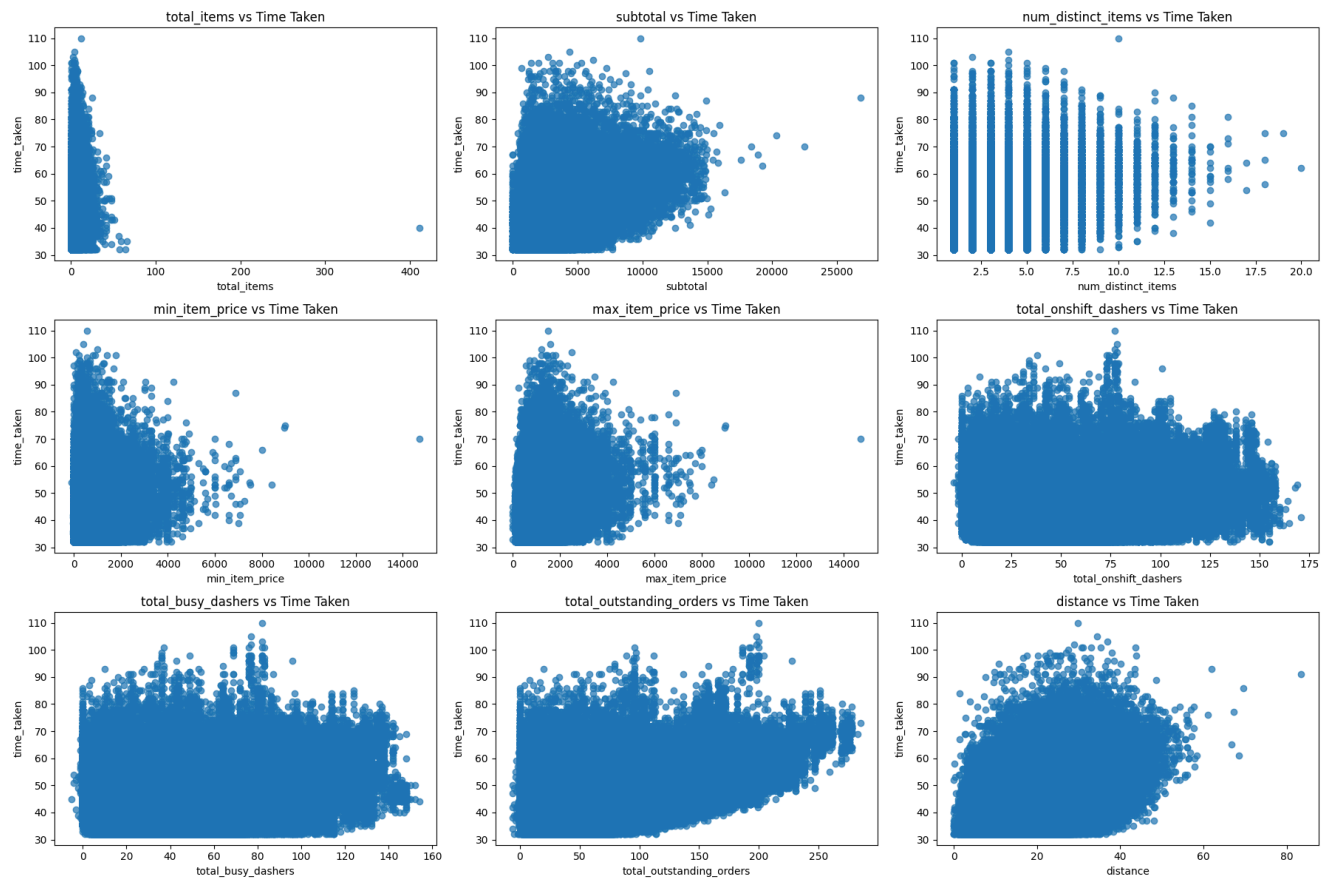
include the number of items, the order subtotal, and the specific time of day (e.g., 4 PM and 5 AM). This suggests that order complexity and operational conditions at specific hours are critical.
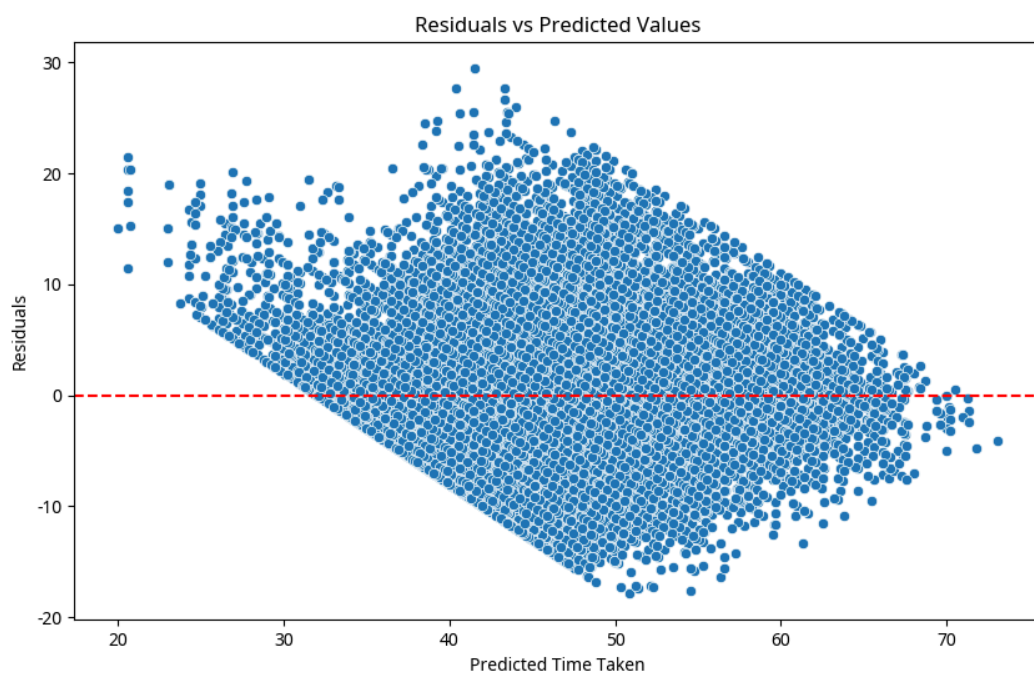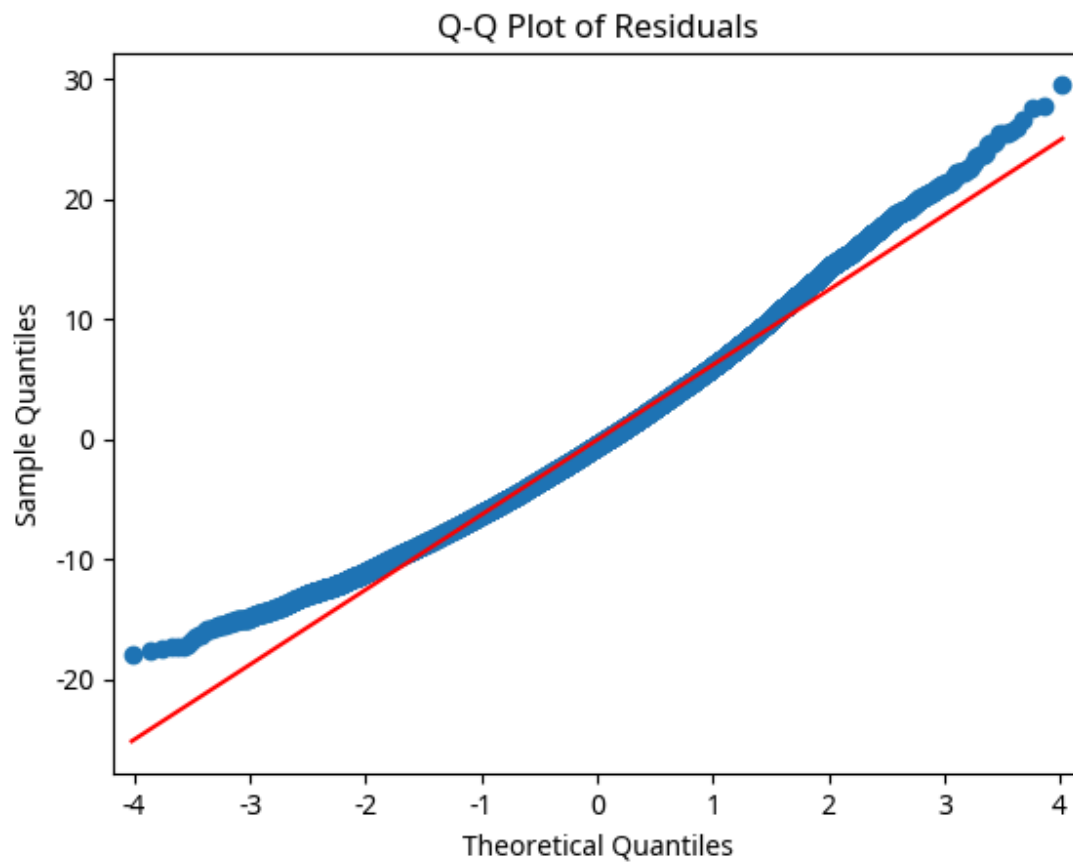
- Operational Optimization: By understanding that certain hours are associated with longer delivery times, the business can proactively adjust staffing levels of delivery partners to meet demand and mitigate delays.
- Customer Experience: While the model is not perfectly predictive, it provides a baseline for more realistic delivery time estimates, which can be used to manage customer expectations and improve satisfaction.

## 8. Visualizations

Correlation Matrix of Features

| total_items vs Time Taken | subtotal vs Time Taken | num_distinct_items vs Time Taken |
| min_item_price vs Time Taken | max_item_price vs Time Taken | total_onshift_dashers vs Time Taken |
| total_busy_dashers vs Time Taken | total_outstanding_orders vs Time Taken | distance vs Time Taken |

Histogram of Residuals

Q-Q Plot of Residuals


Residuals vs Predicted Values

## Time Taken Distribution by Hour

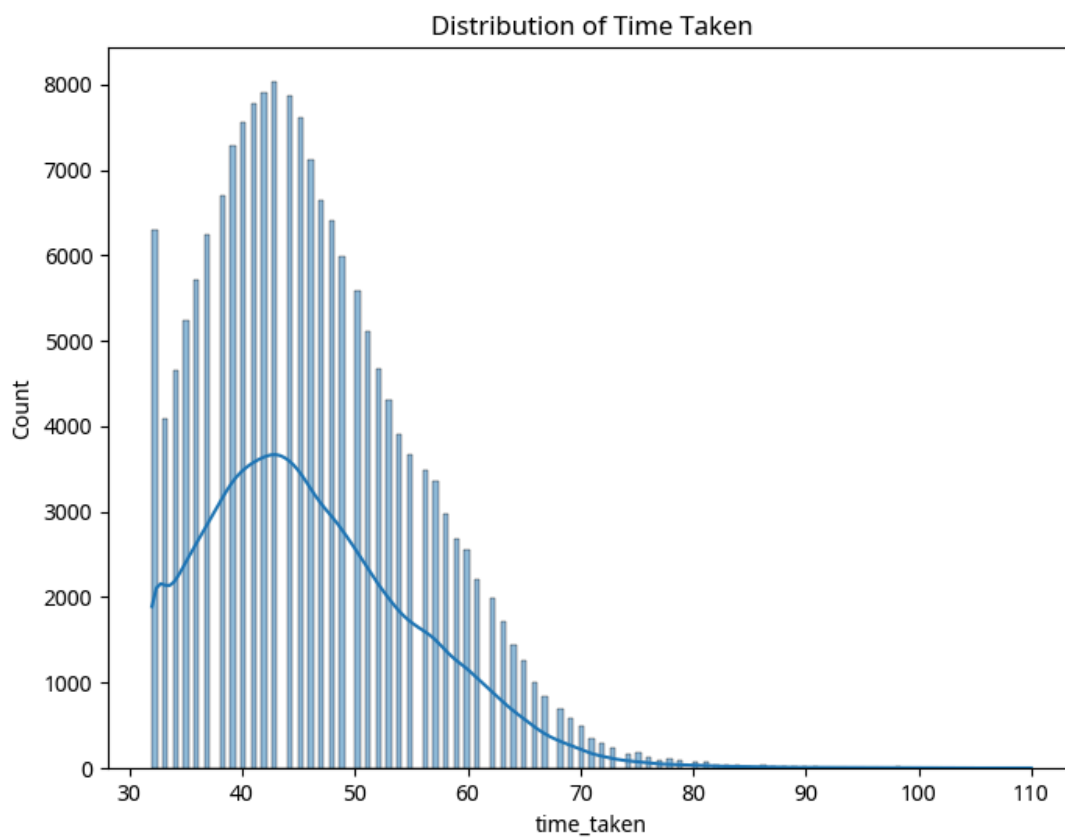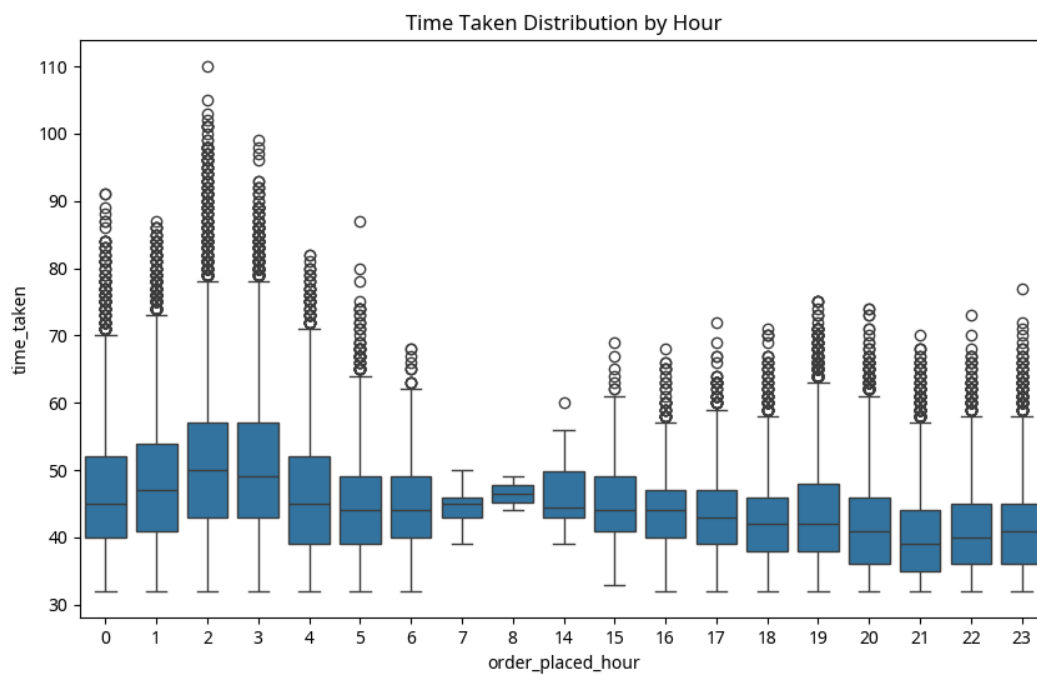

## Distribution of Time Taken

## 8. Conclusion and Future Work

This project successfully developed a linear regression model that provides valuable, interpretable insights into the factors affecting order delivery time. The model serves as a strong foundation for a more sophisticated prediction system.

Future work could focus on:

- Incorporating External Data: Enhancing the model by adding features like real-time traffic data, weather conditions, and local events.
- Advanced Modeling: Experimenting with more complex, non-linear models (e.g., Gradient Boosting, Random Forest) that may capture more intricate patterns in the data and improve the $R^2$ score.

## 9. Conclusion

This project successfully developed a linear regression model to predict order delivery times. While the model's predictive power is moderate, it provides valuable insights into the factors influencing delivery duration. The findings from this analysis can be used to inform business decisions, improve operational efficiency, and enhance the customer experience.