# Station Dwell Time Prediction: A Machine Learning Approach

Umar Aslam

December 28, 2024

**Abstract**

This technical report presents a machine learning-based approach for predicting station dwell times in public transportation systems. We develop a comprehensive model that considers multiple factors including passenger volume, weather conditions, station complexity, and temporal patterns. The model employs ensemble learning techniques, specifically Random Forest and Gradient Boosting, to achieve robust predictions. Our results demonstrate strong predictive performance with an $R^2$ score of approximately 0.84, indicating the model's effectiveness in capturing complex relationships between various factors affecting dwell times.

## 1 Introduction

Station dwell time, defined as the time a transit vehicle spends at a station for passenger boarding and alighting, is a critical factor in public transportation system performance. Accurate prediction of dwell times can significantly improve schedule reliability and overall system efficiency.

## 2 Methodology

### 2.1 Problem Formulation

The dwell time prediction problem can be formulated as a supervised regression task. Given a set of features $X = \{x_1, x_2, ..., x_n\}$ and corresponding dwell times $Y = \{y_1, y_2, ..., y_n\}$, we aim to learn a function $f$ such that:

$$\hat{y} = f(X) + \epsilon \tag{1}$$

where $\hat{y}$ represents the predicted dwell time and $\epsilon$ is the error term.

### 2.2 Feature Engineering

The model incorporates several key features:

- Passenger Volume ($x_{pv}$): Number of passengers per hour

- Platform Length ($x_{pl}$): Platform length in meters

- Peak Hour Indicator ($x_{ph}$): Binary variable (0 or 1)

- Weather Condition ($x_w$): Categorical variable (0: Good, 1: Rain, 2: Snow)

- Station Complexity ($x_{sc}$): Ordinal variable (1: Simple, 2: Medium, 3: Complex)

- Temporal Features: Hour ($x_h$), Month ($x_m$), Day of Week ($x_d$)

- Special Event Indicator ($x_{se}$): Binary variable (0 or 1)

## 2.3 Base Dwell Time Model

The base dwell time calculation follows:

$$t_{base} = t_0 + \alpha x_{pv} + \beta x_{ph} + \gamma x_w + \delta x_{se} + \epsilon x_{sc} + \eta \tag{2}$$

where:

- $t_0$ is the minimum dwell time (60 seconds)

- $\alpha$ is the passenger volume coefficient (0.1)

- $\beta$ is the peak hour impact (30 seconds)

- $\gamma$ is the weather condition impact (15 seconds per level)

- $\delta$ is the special event impact (45 seconds)

- $\epsilon$ is the station complexity impact (10 seconds per level)

- $\eta$ is random noise $\sim \mathcal{N}(0, 10)$

## 2.4 Ensemble Learning Approach

We implement two ensemble learning methods:

### 2.4.1 Random Forest

The Random Forest model uses $K$ decision trees, with each tree's prediction denoted as $h_k(X)$. The final prediction is:

$$\hat{y}_{RF} = \frac{1}{K} \sum_{k=1}^{K} h_k(X) \tag{3}$$

### 2.4.2 Gradient Boosting

The Gradient Boosting model builds trees sequentially, with each tree fitting the residuals of previous trees:

$$\hat{y}_{GB} = \sum_{m=1}^{M} \gamma_m h_m(X) \tag{4}$$

where $\gamma_m$ is the learning rate and $h_m$ is the $m$-th tree.

# 3 Anomaly Detection

We implement statistical anomaly detection using the residual-based approach:

$$A(y, \hat{y}) = \begin{cases} 1 & \text{if } |y - \hat{y}| > \lambda \sigma_\epsilon \\ 0 & \text{otherwise} \end{cases} \tag{5}$$

where $\sigma_\epsilon$ is the standard deviation of residuals and $\lambda$ is the threshold (set to 2).

# 4 Results and Analysis

## 4.1 Model Performance

The model achieves the following performance metrics:

- Root Mean Square Error (RMSE): Approximately 45 seconds

- $R^2$ Score: 0.84

- Mean Absolute Error (MAE): Varies between 30-40 seconds

## 4.2 Feature Importance

Feature importance analysis reveals:

$$I(x_i) = \frac{\sum_{t=1}^{T} n_t I_t(x_i)}{\sum_{t=1}^{T} n_t} \tag{6}$$

where $I_t(x_i)$ is the importance of feature $x_i$ in tree $t$, and $n_t$ is the number of samples reaching a non-leaf node in tree $t$.

Key findings:

- Passenger volume accounts for approximately 48

- Peak hour status contributes around 29

- Weather conditions influence about 17

- Platform length and day of week have minor impacts

# 5 Threshold Alert System

The system generates alerts based on:

$$Alert(x) = \begin{cases} \text{High Volume} & \text{if } x_{pv} > Q_{95}(x_{pv}) \\ \text{Weather Impact} & \text{if } x_w > 0 \\ \text{Special Event} & \text{if } x_{se} = 1 \end{cases} \tag{7}$$

where $Q_{95}(x_{pv})$ represents the 95th percentile of passenger volume.

# 6 Limitations and Future Work

Current limitations include:

- Reliance on synthetic data

- Assumption of independent events

- Limited consideration of infrastructure constraints

Future improvements could include:

- Integration of real-time passenger flow data

- Incorporation of network effects

- Development of platform-specific models

# 7 Conclusion

The developed model demonstrates strong potential for practical application in transit operations. Its ability to capture complex relationships between various factors affecting dwell times makes it a valuable tool for transit planning and real-time operations management.