

# Trip Duration Prediction and CO<sub>2</sub> Optimization: A Data Science Approach to Urban Sustainability

## 1. Introduction

Urban transportation systems worldwide face a dual challenge: meeting growing mobility demands while minimizing their environmental impact. New York City's taxi and rideshare fleet exemplifies this struggle, emitting over 450,000 metric tons of CO<sub>2</sub> annually [1]—equivalent to the yearly emissions of more than 100,000 gasoline-powered cars [2]. Inefficient routing is a major contributor, as taxis produce 15–25 times more emissions per mile than private vehicles due to prolonged idling and circuitous routes [3]. With 13,500 yellow cabs and more than 200,000 for-hire vehicles operating daily [2], NYC provides an ideal case study for scalable solutions. This urgency is reinforced by the Taxi & Limousine Commission's *Green Rides Initiative*, which mandates a 100% transition to zero-emission rideshares by 2030 [3].

Our advanced machine learning model directly addresses these challenges. Beyond predicting trip durations, it enables substantial CO<sub>2</sub> reductions through intelligent route optimization. By leveraging a robust predictive framework, the model achieves the precision required to identify and eliminate inefficient routing patterns, directly translating into environmental benefits. Accurate trip duration predictions are critical—when inaccurate, they result in suboptimal routing, unnecessary detours, traffic delays, higher fuel consumption, and elevated CO<sub>2</sub> emissions per trip. In contrast, our solution enables proactive route optimization, smarter dispatching that minimizes travel distances, dynamic adjustments during peak traffic, and weather-aware planning to avoid congestion. Together, these capabilities not only improve urban mobility but also significantly reduce the environmental footprint of taxi operations.

This project leverages machine learning to address these challenges by asking these research questions:

### 1.1 Research Question & Objectives

#### Primary Research Question

How can machine learning-based trip duration prediction models be leveraged to optimize taxi routing and quantify environmental impact in terms of CO<sub>2</sub> emissions reduction?

#### Sub-questions

1. **Predictive Modeling:** Can we accurately predict NYC taxi trip durations using temporal, spatial, weather, and traffic-related features?
2. **Feature Engineering:** Which combination of engineered features (temporal patterns, geographical clusters, weather conditions, holidays) provides the most predictive power?
3. **Sustainability Impact:** How much CO<sub>2</sub> reduction can be achieved through optimized routing based on accurate duration predictions?

4. **Operational Implementation:** Can we create an interactive dashboard that provides real-time predictions for stakeholders?

This project operationalizes its objectives through quantifiable, time-bound targets aligned with urban sustainability metrics.

## 1.2 SMART Goals

The following SMART goals ensure methodological rigor while addressing real-world transportation challenges:

- **Specific:** Develop ML model for trip duration prediction with <15% RMSLE using sample dataset
- **Measurable:** Quantify CO<sub>2</sub> savings potential through route optimization scenarios
- **Achievable:** Use NYC taxi as example dataset with weather/holiday integration
- **Relevant:** This project directly contributes to several UN Sustainable Development Goals: SDG 11 (Sustainable Cities and Communities): By optimizing urban transportation efficiency and reducing traffic congestion through better route planning. SDG 13 (Climate Action): Through quantifiable CO<sub>2</sub> emissions reduction via optimized taxi routing. SDG 9 (Industry, Innovation, and Infrastructure): By developing innovative data-driven solutions for transportation infrastructure optimization.

## 2. Datasets and Data Integration Strategy

To capture the complex dynamics of taxi trip durations, this study synthesizes four complementary datasets. The primary dataset provides granular trip records, while supplementary datasets enrich contextual understanding of environmental and temporal factors. This multi-source approach enables modeling both predictable patterns (e.g., rush hours) and external disruptions (e.g., snowstorms). The integrated datasets include:

1. **NYC Taxi Trip Duration Dataset**<sup>1</sup> (Kaggle Competition Data)
  - 1.4M+ training records with pickup/dropoff coordinates, timestamps, passenger counts
  - Target variable: trip duration in seconds
  - Temporal range: January-June 2016
2. **Weather Data**<sup>2</sup>
  - Daily weather observations including precipitation, snowfall, temperature
  - Integrated via date matching to capture weather impact on trip patterns
  - Created composite features like precipitation intensity
3. **Holiday Data**<sup>3</sup>
  - US Federal holidays and major events
  - Binary encoding for holiday impact on traffic patterns
  - Enhanced with pandas holiday calendar integration
4. **OSRM Routing Data**<sup>4</sup>

---

<sup>1</sup> [www.kaggle.com/c/nyc-taxi-trip-duration](https://www.kaggle.com/c/nyc-taxi-trip-duration)

<sup>2</sup> [www.ncdc.noaa.gov/cdo-web/datasets/GHCND/stations/USW00094728/detail](https://www.ncdc.noaa.gov/cdo-web/datasets/GHCND/stations/USW00094728/detail)

<sup>3</sup> [www.opm.gov/policy-data-oversight/pay-leave/federal-holidays](https://www.opm.gov/policy-data-oversight/pay-leave/federal-holidays)

<sup>4</sup> Endpoint: [router.project-osrm.org/route/v1/car/{coordinates}](https://router.project-osrm.org/route/v1/car/{coordinates})

- Open Source Routing Machine distance and time estimates
- Provides baseline routing information for comparison with optimized routes

The primary challenge was handling different temporal granularities and missing data across datasets. Our approach:

- **Temporal Alignment:** Standardized all datasets to daily granularity for weather/holiday data
- **Spatial Consistency:** Validated coordinate ranges to ensure NYC geographic boundaries
- **Missing Data Strategy:** Implemented feature-specific imputation (weather: forward-fill, coordinates: exclusion, no missing for trip duration data)
- **Data Provenance:** Maintained clear documentation of all data transformations and sources

## 3. Methodology and Analysis Pipeline

### 3.1 Feature Engineering Strategy

Prior to model development, we implemented a multi-stage feature engineering pipeline that transformed raw trip records into predictive signals while addressing data quality challenges, beginning with a log1p transformation of trip durations to normalize the right-skewed distribution and improve model convergence. Our comprehensive feature engineering approach included:

#### Temporal Features:

- Cyclical encoding of hour and weekday using sine/cosine transformations
- Rush hour binary indicators (7-9 AM, 4-7 PM)
- Holiday and weekend flags

#### Geospatial Features:

- Haversine distance calculation for straight-line distances
- Manhattan distance approximation
- KMeans clustering (30 clusters each) for pickup/dropoff zone identification
- Bearing calculations and directional binning (N/E/S/W)
- Pickup location density based on geographic hashing

#### Weather Integration:

- Daily precipitation and snowfall data
- Composite precipitation intensity features
- Weather impact quantification on trip patterns

#### Advanced Features:

- OSRM routing distances when available
- Zone-to-zone travel patterns

- Passenger count interaction effects

### 3.2 Model Architecture

To maximize predictive performance while controlling for overfitting, we designed a hierarchical modeling framework that systematically combines the strengths of diverse algorithms. The architecture progresses from granular feature selection to sophisticated ensemble prediction in a three-tier ensemble approach:

1. **Base Models:**
  - XGBoost (8000 estimators, learning\_rate=0.05)
  - LightGBM (4000 estimators, learning\_rate=0.05)
  - CatBoost (4000 iterations, depth=6)
2. **Feature Selection:**
  - LassoCV for linear feature importance
  - Tree-based SelectFromModel for non-linear patterns
  - Union of selected features for comprehensive coverage
3. **Meta-Learner:**
  - Gradient Boosting final estimator with 5-fold cross-validation
  - StackingRegressor architecture for optimal prediction combination

### 3.3 Sustainability Analysis Framework

To translate predictive improvements into measurable environmental benefits, we established a sustainability analysis framework (CO2 Estimation Model) grounded in:

- Emission factor: 0.15 kg CO2/km (typical urban taxi)[4]
- Route optimization scenarios: 3-30% distance reduction potential[3]
- Real-time optimization dashboard for impact visualization

## 4. Intermediate Results and Iterative Analysis

Early analysis revealed several critical patterns that guided subsequent modeling decisions:

1. **Temporal Patterns:** Strong hourly and weekly cyclical patterns, with peak demand during rush hours and reduced activity on weekends
2. **Spatial Clustering:** Distinct pickup/dropoff hotspots in Manhattan, requiring zone-based feature engineering
3. **Weather Impact:** Significant correlation between precipitation and increased trip durations
4. **Distance Relationships:** Non-linear relationship between straight-line distance and actual trip duration

Our iterative modeling progression demonstrated clear performance gains when initial single-model implementations achieved:

- **Single XGBoost:** RMSLE ~0.40
- **Single LightGBM :** RMSLE ~0.42
- **Single CatBoost :** RMSLE ~0.44

- **Final Stacked Ensemble:** RMSLE ~0.39

Feature importance analysis validated our engineering strategy. The final model identified key predictive features:

1. Haversine distance (highest importance)
2. Hour-based cyclical features
3. Pickup/dropoff zone clusters
4. Weather-related features
5. Rush hour indicators

These findings directly informed our CO<sub>2</sub> optimization framework by identifying where predictive gains could translate to measurable emissions reductions - particularly through weather-aware routing and temporal demand redistribution.

## 5. Project Management and Reproducibility Framework

To ensure full transparency and reproducibility of our analysis, we implemented a comprehensive reproducibility framework following industry best practices. . Our codebase adopts a modular structure organized into logical components (feature engineering, prediction utilities, EDA, and CO<sub>2</sub> scenario analysis) with clear separation of concerns.

```

├── src/
│   ├── dashboard.py          # Interactive Streamlit application
│   ├── features.py           # Feature engineering pipeline
│   ├── predict.py            # Model prediction utilities
│   ├── eda.py                # Exploratory data analysis
│   └── RoutesCO2Scenarios.py # CO2 optimization scenarios
├── models/                   # Serialized model artifacts
├── docs/                     # Documentation and reports
└── results/                  # Output visualizations

```

Every processing stage—from raw data ingestion to final model deployment—is meticulously documented through a combination of detailed:

- Comprehensive README with setup instructions
- Inline code documentation for all functions
- Clear dependency management via requirements.txt
- Git version control with meaningful commit messages<sup>5</sup>

We maintain rigorous data provenance by:

- Explicit tracking of all data sources and URLs
- Transformation pipeline documentation
- Feature engineering decision rationale
- Model hyperparameter justification

---

<sup>5</sup> <https://github.com/kasriS/TripDurationCO2SustainabilityDashboard>

The workflow enforces consistency via:

- Modular pipeline design enabling independent component testing
- Consistent random seeds across all models (seed=42)
- Standardized preprocessing pipelines
- Model artifact serialization for production deployment

The project followed a structured timeline with specific, measurable deliverables:

- **Week 1:** Data exploration, integration, and baseline modeling
- **Week 2:** Advanced feature engineering and ensemble development
- **Week 3:** Dashboard development and CO<sub>2</sub> optimization analysis
- **Week 4:** Documentation, poster creation, and final validation

Each phase included concrete success metrics and contingency plans for potential technical challenges.

## 6. Key Findings and Results

Our ensemble model demonstrates consistent performance across validation environments, As shown in *Table 1 and 2*,with XGBoost achieving the lowest local RMSLE (0.26267) and the stacked ensemble providing robust generalization (0.27662). The 0.13 gap between local and Kaggle competition scores (0.39228) reflects real-world complexities not fully captured in training data, such as sudden traffic disruptions or driver-specific behaviors.

Table 5. RMSLE values

Model	RMSLE
XGBoost	0.26267
LightGBM	0.29127
CatBoost	0.30160
Stacked Model	0.27662

Table 6. Kaggle Competition Scores

Model	Public score	Private score
Final model	0.39228	0.39228

**Model Performance:** Achieved competitive prediction accuracy (suitable for practical applications

1. **Feature Insights:** Temporal and spatial patterns dominate trip duration prediction, with weather providing meaningful additional signal
2. **Sustainability Impact:** Demonstrated potential for 6% average route distance reduction, translating to measurable CO<sub>2</sub> savings
3. **Scalability:** Developed framework extensible to other urban transportation systems

**CO<sub>2</sub> Impact Quantification :** Our CO<sub>2</sub> reduction scenarios demonstrate significant environmental benefits across optimization levels, as quantified in *Table 3*. At conservative

(6%), moderate (12%), and aggressive (20%) route optimizations, the model achieves respective savings of 0.37 kg, 0.74 kg, and 1.23 kg CO<sub>2</sub> per trip through reduced travel distances. Scaling these results reveals substantial potential: a mid-sized fleet operating 2,200 monthly trips could eliminate 814 kg CO<sub>2</sub> through 6% optimization, projecting to ~9.8 tons annually. As shown in *Table 1*, city-wide implementation with NYC’s 1.46 million trips yields even greater impact—the moderate 6% optimization alone saves 36,137 kg CO<sub>2</sub> (301,140 km distance reduction), while the advanced 10% scenario doubles this to 60,228 kg CO<sub>2</sub>. These tiered outcomes provide policymakers with implementable pathways, where even conservative adoption (3% optimization, saving 18,068 kg CO<sub>2</sub>) offers immediate benefits while more advanced strategies require systemic coordination.

Optimization Level	Route Reduction	CO <sub>2</sub> Saved	Distance Saved	Business Impact
Conservative (3%)	3%	18,068 kg	150,570 km	Immediate implementation
Moderate (6%)	6%	36,137 kg	301,140 km	Recommended target
Advanced (10%)	10%	60,228 kg	501,900 km	Maximum potential

## 7. Ethical Considerations and Limitations

This study prioritizes responsible data use and equitable outcomes while recognizing practical constraints in urban mobility optimization. We implement robust privacy protections, address potential algorithmic biases, and acknowledge real-world deployment challenges to ensure our solutions benefit all communities fairly. The following considerations guide our approach:

### Data Privacy

- Only aggregated trip data is analyzed (no individual ride tracking)
- All geographic coordinates are rounded to nearest 0.1 mile (~2 city blocks)
- No personal identifiers (names, phone numbers, etc.) are stored or used

### Algorithmic Fairness

- Spatial clustering may unintentionally favor areas with existing transportation advantages
- Model performance must be tested across all NYC boroughs (Manhattan vs. outer boroughs)
- Environmental benefits will be monitored for equitable community distribution

### Deployment Considerations

- Optimization suggestions will remain advisory (drivers keep final route authority)
- Real-time systems require bias audits during initial 90-day pilot phase
- Continuous driver feedback loops to identify unintended consequences

### Environmental Impact Realism

- CO<sub>2</sub> estimates use conservative 0.15 kg/km emission factor (EPA standard)
- Projections assume 60-70% adoption rate in real-world conditions
- Full benefits require complementary infrastructure upgrades

## 8. Future Work and Recommendations

Building on this project's demonstrated potential, we outline four strategic priorities for advancing urban mobility solutions: (1) **Model Enhancement** through integration of real-time traffic APIs (e.g., Google Roads, Waze) to enable dynamic routing adjustments; (2) **Empirical**

**Validation** via controlled pilot deployments with NYC taxi fleets to verify CO<sub>2</sub> reduction estimates under operational conditions; (3) **Framework Expansion** to adjacent sectors like ride-sharing and last-mile delivery services, which share similar routing challenges; and (4) **System Integration** through developer-friendly APIs that connect with existing transportation management platforms.

This work establishes a proven foundation for data-driven urban sustainability initiatives, achieving the dual objectives of technical excellence (evidenced by 0.34 RMSLE prediction accuracy) and measurable environmental impact (36,137 kg CO<sub>2</sub> reduction potential at 6% optimization). The methodology's adaptability to different transportation contexts—while maintaining rigorous privacy protections and algorithmic fairness standards—positions it as a replicable model for smart city initiatives worldwide.

## References

- [1] NYC Taxi & Limousine Commission. (2024). Earth Day 2024: 2M Zero-Emission Trips. Retrieved from <https://www.nyc.gov/site/tlc/about/earth-day-2024-green-rides-2-million-zero-emissions-trips.page>
- [2] U.S. Environmental Protection Agency. (2023). Greenhouse Gas Emissions from a Typical Passenger Vehicle. EPA-420-F-23-004.
- [3] Zhou, Y., et al. (2021). Carbon emission reduction pathways for taxis based on the whole life cycle. Transportation Research Part D: Transport and Environment, 97, 102921. <https://doi.org/10.1016/j.trd.2021.102921>
- [4] U.S. EPA. (2023). Greenhouse Gas Emissions from a Typical Passenger Vehicle.