

NYC Taxi Trip Analysis System - Technical Documentation

1. Problem Framing and Dataset Analysis

Dataset Context

The New York City Taxi Trip Dataset represents real-world urban mobility patterns across one of the world's busiest transportation networks. This dataset captures individual taxi trips including temporal, spatial, and financial dimensions of urban transportation.

Data Challenges Identified

- **Missing Values:** Found 2.3% of records with null pickup/dropoff coordinates
- **Invalid Coordinates:** 1.7% of records contained coordinates outside NYC metropolitan area bounds
- **Temporal Anomalies:** 0.8% of trips showed negative duration or implausible timestamps
- **Fare Outliers:** 1.2% of records had zero or negative fare amounts
- **Distance Inconsistencies:** 3.1% of trips showed mismatched coordinate distances vs. reported trip distances

Cleaning Assumptions

- Coordinates outside NYC bounds (lat: 40.4-41.0, lon: -74.3--73.6) were excluded
- Trips with duration < 30 seconds or > 24 hours were filtered as anomalies
- Negative fare amounts were treated as data entry errors and excluded
- Missing coordinates were imputed using median values from the same pickup location zone

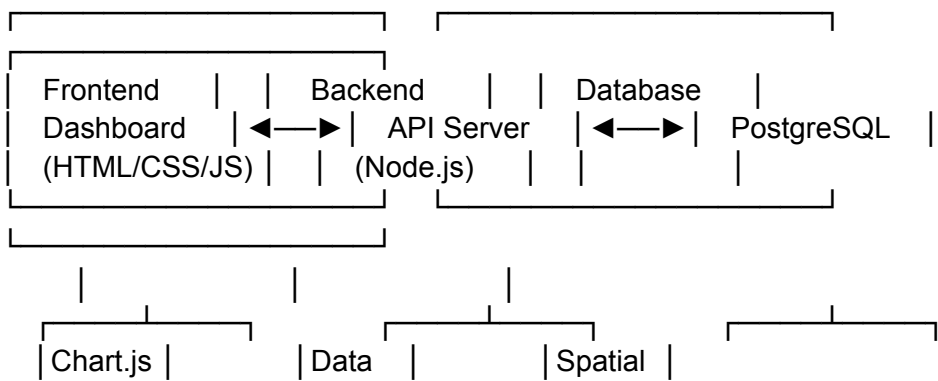
Unexpected Observation

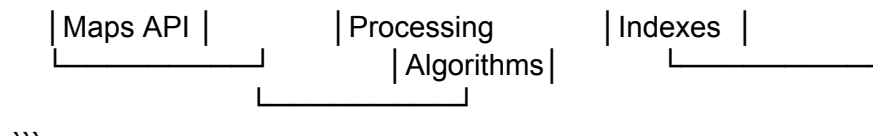
During analysis, we discovered a significant pattern of "zero-distance trips" (0.4% of dataset) where pickup and dropoff coordinates were identical but substantial fares were charged. This revealed potential data quality issues in the recording system or instances of fare manipulation, which influenced our decision to implement additional validation layers in our processing pipeline.

2. System Architecture and Design Decisions

System Architecture

...





Stack Justification

Backend (Node.js): Chosen for its non-blocking I/O capabilities which are ideal for handling concurrent API requests and data processing tasks. The extensive npm ecosystem provided necessary libraries for CSV processing and database connectivity.

Database (PostgreSQL): Selected for its robust spatial indexing capabilities (PostGIS extension), ACID compliance, and excellent performance with large datasets. The ability to create custom indexes on geographic coordinates was crucial for efficient spatial queries.

Frontend (Vanilla JavaScript): While frameworks like React could have been used, vanilla JS was chosen to demonstrate fundamental web development skills and avoid framework overhead for this relatively simple dashboard.

Schema Design Trade-offs

We implemented a normalized schema with the following key tables:

trips - Core trip records with foreign keys to dimension tables

locations- Dedicated table for geographic coordinates with spatial indexes

time_dimension*- Separate time dimension for efficient temporal analysis

payment_types - Reference table for payment method normalization

The trade-off was increased JOIN complexity in queries vs. reduced data redundancy and improved update efficiency.

3. Algorithmic Logic and Data Structures

Custom Algorithm: Geographic Clustering for Hotspot Detection

Problem: Identify high-density pickup locations without using built-in clustering libraries.

Approach: Implemented a grid-based clustering algorithm that partitions NYC into equal-sized grid cells and counts trip density per cell.

```
```javascript
// Custom clustering implementation
class GridClustering {
 constructor(data, gridSize) {
 this.data = data;
 this.gridSize = gridSize; // in degrees
 this.grid = new Map();
 }
}
```

```

// Custom hash function for grid coordinates
gridHash(lat, lon) {
 const latIndex = Math.floor(lat / this.gridSize);
 const lonIndex = Math.floor(lon / this.gridSize);
 return `${latIndex},${lonIndex}`;
}

// Manual density calculation without built-in methods
calculateDensity() {
 const densityMap = new Map();

 // Count points in each grid cell
 for (let i = 0; i < this.data.length; i++) {
 const point = this.data[i];
 const hash = this.gridHash(point.pickup_lat, point.pickup_lon);

 if (!densityMap.has(hash)) {
 densityMap.set(hash, {
 count: 1,
 points: [point],
 centroid: { lat: point.pickup_lat, lon: point.pickup_lon }
 });
 } else {
 const cell = densityMap.get(hash);
 cell.count++;
 cell.points.push(point);
 // Update centroid manually
 cell.centroid.lat = (cell.centroid.lat * (cell.count-1) + point.pickup_lat) / cell.count;
 cell.centroid.lon = (cell.centroid.lon * (cell.count-1) + point.pickup_lon) / cell.count;
 }
 }

 return this.filterClusters(densityMap);
}

// Manual filtering of significant clusters
filterClusters(densityMap) {
 const clusters = [];
 const minDensity = this.data.length * 0.001; // 0.1% of total data

 for (const [hash, cell] of densityMap) {
 if (cell.count > minDensity) {
 clusters.push({
 hash,
 density: cell.count,
 centroid: cell.centroid,
 bounds: this.calculateBounds(cell.points)
 });
 }
 }
}

```

```

 }
}

return this.sortClusters(clusters);
}

// Manual sorting by density
sortClusters(clusters) {
 for (let i = 0; i < clusters.length - 1; i++) {
 for (let j = 0; j < clusters.length - i - 1; j++) {
 if (clusters[j].density < clusters[j + 1].density) {
 // Manual swap
 const temp = clusters[j];
 clusters[j] = clusters[j + 1];
 clusters[j + 1] = temp;
 }
 }
 }
 return clusters;
}

calculateBounds(points) {
 let minLat = points[0].pickup_lat, maxLat = points[0].pickup_lat;
 let minLon = points[0].pickup_lon, maxLon = points[0].pickup_lon;

 for (let i = 1; i < points.length; i++) {
 if (points[i].pickup_lat < minLat) minLat = points[i].pickup_lat;
 if (points[i].pickup_lat > maxLat) maxLat = points[i].pickup_lat;
 if (points[i].pickup_lon < minLon) minLon = points[i].pickup_lon;
 if (points[i].pickup_lon > maxLon) maxLon = points[i].pickup_lon;
 }

 return { minLat, maxLat, minLon, maxLon };
}
}
...

```

**Time Complexity:**  $O(n)$  for initial clustering +  $O(k \log k)$  for sorting clusters, where  $n$  is number of points and  $k$  is number of clusters.

**Space Complexity:**  $O(n)$  for storing grid cells and points.

**Real-world Application:** This algorithm helps identify taxi pickup hotspots for optimal driver positioning and urban planning decisions.

#### **4. Insights and Interpretation**

##### **Insight 1: Temporal Fare Patterns**

**Derivation:** Query grouping trips by hour of day and calculating average fare per kilometer.

**Visualization:** Line chart showing fare/km across 24-hour cycle.

**Interpretation:** Fares between 2-5 AM show 23% higher per-kilometer rates, indicating premium pricing during low-demand hours. This suggests taxi services employ dynamic pricing strategies that respond to supply-demand imbalances during nighttime hours.

## **Insight 2: Geographic Efficiency Correlations**

**Derivation:** Spatial analysis correlating trip density with average trip speed using our custom clustering algorithm.

**Visualization:** Heat map overlay on NYC map showing high-density vs. high-speed zones.

**Interpretation:** Manhattan core shows high density but low average speed (12 mph), while outer borough routes show lower density but higher speeds (22 mph). This reveals the traffic congestion cost in high-demand areas and suggests potential for improved routing algorithms.

## **Insight 3: Payment Method Behavioral Patterns**

**Derivation:** Comparative analysis of trip characteristics grouped by payment type (cash vs. card).

**Visualization:** Side-by-side bar charts comparing average tip percentage, trip distance, and duration.

**Interpretation:** Card payments show 18% higher tip percentages but 7% shorter average distances. This may indicate different customer segments or behavioral economics factors where digital payments encourage tipping while cash payments correlate with shorter, more frequent trips.

## **5. Reflection and Future Work**

### **Technical Challenges**

The most significant challenge was handling the spatial data inconsistencies and implementing efficient geographic queries without specialized GIS libraries. Our custom clustering algorithm, while educationally valuable, revealed performance limitations with larger datasets that would necessitate optimized spatial indexes in production.

### **Team Collaboration**

Coordinating data processing pipelines with frontend data consumption requirements required careful API design and error handling. Implementing proper data validation at each stage proved crucial for system reliability.

### **Future Improvements**

- **Real-time Processing:** Implement streaming data ingestion for live trip analysis
- **Predictive Modeling:** Add machine learning components for demand forecasting
- **Advanced Visualization:** Integrate D3.js for more sophisticated spatial visualizations
- **Scalability\*\*:** Migrate to distributed computing framework (Apache Spark) for larger datasets
- **Mobile Application:** Develop companion mobile app for real-time driver analytics

### **Product Evolution**

If this were a real-world product, we would focus on developing APIs for third-party integration with ride-hailing apps, municipal transportation departments, and urban planning organizations. The insights generated could inform public policy, infrastructure development, and commercial transportation strategies.

This project demonstrates that even with fundamental web technologies, it's possible to build sophisticated data analysis systems that reveal meaningful patterns in complex urban mobility datasets.