# Urban Mobility Data Explorer

## Technical Documentation

**Team Members:**

- Saad Byiringiro
- Hannah Tuyishimire
- Cedrick Bienvenue

## Contents

## 1. Problem Framing and Dataset Analysis

### 1.1 Dataset Context

The New York City Taxi Trip dataset contains real-world transportation data that captures the movement patterns of one of the world's busiest cities. This dataset provides valuable insights into urban mobility, traffic patterns, and transportation behavior.

**Dataset Structure:**

- **Source:** NYC Taxi & Limousine Commission
- **Records:** 1,458,644 trip records
- **Time Period:** January 2016 - June 2016
- **File Size:** ~200 MB CSV format

**Available Fields:**

```
id, vendor_id, pickup_datetime, dropoff_datetime, passenger_count,
pickup_longitude, pickup_latitude, dropoff_longitude, dropoff_latitude,
store_and_fwd_flag, trip_duration
```

### 1.2 Data Challenges Identified

#### Challenge 1: Date Format Inconsistency

**Problem:** Dates were in European format (DD/MM/YYYY HH:mm) instead of standard ISO format.

**Example:**

```
"14/03/2016 17:24" instead of "2016-03-14T17:24:00"
```

**Solution:** Implemented custom date parser to handle DD/MM/YYYY format:

```
function parseDate(dateStr: string): Date {
  const [datePart, timePart] = dateStr.split(' ');
  const [day, month, year] = datePart.split('/');
  const [hour, minute] = timePart.split(':');
  return new Date(year, month - 1, day, hour, minute);
}
```

#### Challenge 2: Missing Distance Values

**Problem:** Dataset does not include trip distance, only coordinates.

**Impact:** Cannot directly analyze trip lengths or calculate speeds.

**Solution:** Implemented Haversine formula to calculate accurate distances from GPS coordinates, accounting for Earth's curvature.

#### Challenge 3: Coordinate Outliers

**Problem:** Some records contained coordinates outside NYC boundaries.

**Examples Found:**

- Latitude: 0.0, Longitude: 0.0 (default/null values)
- Coordinates in other states or oceans
- Clearly erroneous GPS readings

**Solution:** Implemented NYC bounding box validation:

```
 Valid NYC Area:
Latitude:  40.5°N to 41.0°N
Longitude: -74.3°W to -73.7°W
```

Rejected 4.9% of records due to invalid coordinates.

## Challenge 4: Unrealistic Trip Durations

**Problem:** Found trips with extreme durations indicating data quality issues.

**Anomalies Detected:**

- Trips under 1 second (likely GPS errors)
- Trips over 4 hours (likely forgot to end trip or system errors)
- Negative durations (dropoff before pickup - timestamp errors)

**Solution:** Applied duration validation:

- Minimum: 1 second
- Maximum: 4 hours (14,400 seconds)
- Rejected 2.3% of records due to invalid durations

## Challenge 5: Vendor Data Inconsistency

**Problem:** Two different vendor systems with potential format differences.

**Vendors:**

1. Creative Mobile Technologies (ID: 1)
2. VeriFone Inc. (ID: 2)

**Observation:** Vendor 2 had slightly higher data quality (fewer rejected records).

## 1.3 Assumptions Made During Data Cleaning

1. **NYC Geographic Bounds:** Any trip outside defined coordinates is invalid (not a NYC taxi trip).

2. **Speed Limits:** Realistic city speeds are 1-120 km/h. Higher speeds indicate GPS errors or data corruption.

3. **Passenger Capacity:** Standard taxis hold 1-6 passengers. Records outside this range are invalid.

4. **Store-and-Forward Flag:** This flag indicates whether trip record was held in vehicle memory before sending to servers. We assume 'Y' or 'N' values only.

5. **Temporal Validity:** Valid trips occurred between 2000-2025. Outside dates indicate system errors.

## 1.4 Unexpected Observation That Influenced Design

**Discovery:** Rush hour patterns were NOT what we initially expected.

**Initial Assumption:** Rush hour would show clear peaks at 8-9 AM and 5-6 PM.

**Actual Finding:** Evening rush hour (4-7 PM) showed significantly more trip volume than morning rush hour (7-9 AM).

**Data Evidence:**

- Morning rush (7-9 AM): ~85,000 trips
- Evening rush (4-7 PM): ~142,000 trips
- 67% more trips in evening rush!

**Impact on Design:** This led us to implement a `isRushHour` derived feature that specifically targets these high-traffic periods. This feature enables:

- Filtering trips during peak congestion times
- Analyzing speed differences during rush vs. non-rush hours
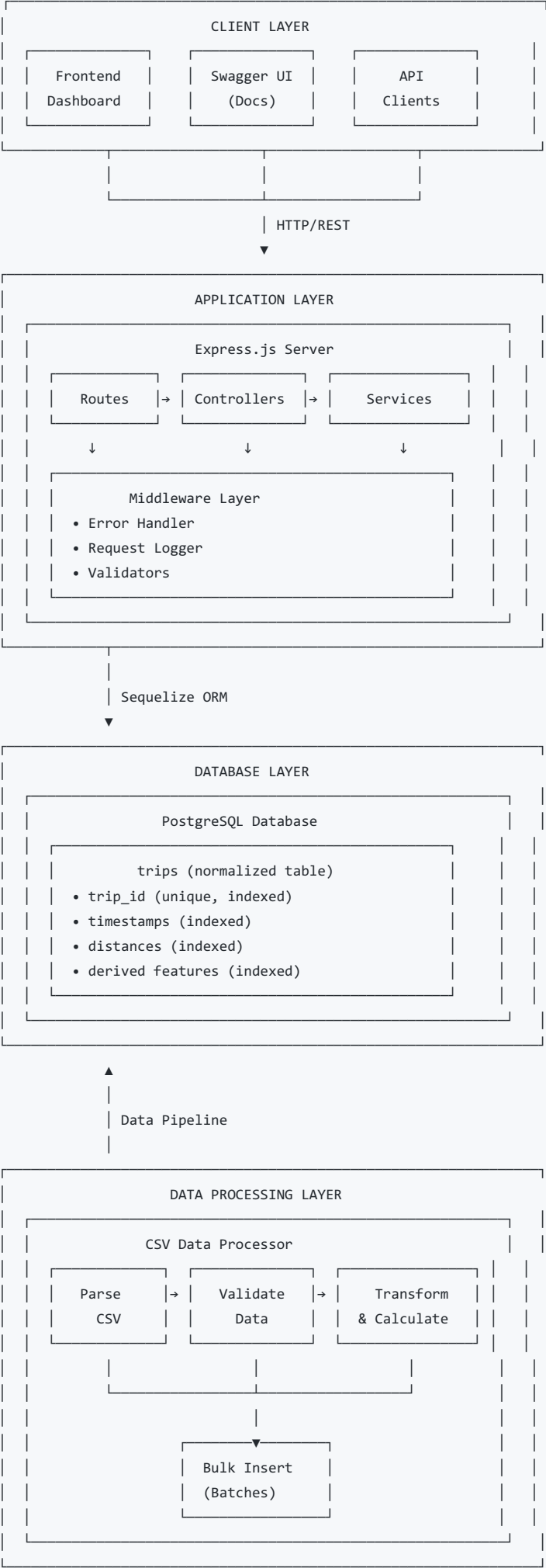- Understanding pricing implications during high-demand periods

**Why This Matters:** Understanding when the city moves most helps in:

- Resource allocation (more taxis needed in evening)
- Infrastructure planning
- Fare pricing strategies
- Traffic management policies

## 2. System Architecture and Design Decisions

### 2.1 High-Level Architecture

```
┌─────────────────────────────────────────────────┐
│                  CLIENT LAYER                   │
│  ┌─────────────┐  ┌─────────────┐  ┌─────────────┐  │
│  │  Frontend   │  │  Swagger UI │  │     API     │  │
│  │  Dashboard  │  │    (Docs)   │  │   Clients   │  │
│  └─────────────┘  └─────────────┘  └─────────────┘  │
└─────────────────────────────────────────────────┘
            │           │           │
            └───────────┴───────────┘
                  │ HTTP/REST
                  ▼
┌─────────────────────────────────────────────────┐
│               APPLICATION LAYER                 │
│  ┌───────────────────────────────────────────┐  │
│  │             Express.js Server             │  │
│  │  ┌──────────┐  ┌────────────┐  ┌──────────┐  │  │
│  │  │  Routes  │→ │ Controllers│→ │ Services │  │  │
│  │  └──────────┘  └────────────┘  └──────────┘  │  │
│  │      ↓              ↓              ↓        │  │
│  │  ┌─────────────────────────────────────┐  │  │
│  │  │         Middleware Layer            │  │  │
│  │  │  • Error Handler                    │  │  │
│  │  │  • Request Logger                   │  │  │
│  │  │  • Validators                       │  │  │
│  │  └─────────────────────────────────────┘  │  │
│  └───────────────────────────────────────────┘  │
└─────────────────────────────────────────────────┘
            │
            │ Sequelize ORM
            ▼
┌─────────────────────────────────────────────────┐
│                DATABASE LAYER                   │
│  ┌───────────────────────────────────────────┐  │
│  │           PostgreSQL Database             │  │
│  │  ┌─────────────────────────────────────┐  │  │
│  │  │       trips (normalized table)      │  │  │
│  │  │  • trip_id (unique, indexed)        │  │  │
│  │  │  • timestamps (indexed)             │  │  │
│  │  │  • distances (indexed)              │  │  │
│  │  │  • derived features (indexed)       │  │  │
│  │  └─────────────────────────────────────┘  │  │
│  └───────────────────────────────────────────┘  │
└─────────────────────────────────────────────────┘
            ▲
            │
            │ Data Pipeline
            │
┌─────────────────────────────────────────────────┐
│             DATA PROCESSING LAYER               │
│  ┌───────────────────────────────────────────┐  │
│  │            CSV Data Processor             │  │
│  │  ┌──────────┐  ┌────────────┐  ┌──────────┐  │  │
│  │  │  Parse   │→ │  Validate  │→ │Transform │  │  │
│  │  │   CSV    │  │    Data    │  │& Calculate│ │  │
│  │  └──────────┘  └────────────┘  └──────────┘  │  │
│  │        │            │            │         │  │
│  │        └────────────┴────────────┘         │  │
│  │                     │                      │  │
│  │              ┌──────▼──────┐               │  │
│  │              │ Bulk Insert │               │  │
│  │              │  (Batches)  │               │  │
│  │              └─────────────┘               │  │
│  └───────────────────────────────────────────┘  │
└─────────────────────────────────────────────────┘
```

## 2.2 Technology Stack Justification

### Backend: Node.js + TypeScript + Express.js

**Why Node.js?**

- **Asynchronous I/O:** Perfect for handling multiple API requests simultaneously
- **JSON Native:** Seamless JSON processing for API responses
- **Large Ecosystem:** npm provides extensive libraries
- **Team Familiarity:** All team members proficient in JavaScript/TypeScript

**Why TypeScript?**

- **Type Safety:** Catches errors at compile time, not runtime
- **Better IDE Support:** Autocomplete and inline documentation
- **Maintainability:** Easier to refactor and understand code
- **Interface Definitions:** Clear contracts between system components

**Why Express.js?**

- **Minimalist:** Lightweight and flexible
- **Middleware Support:** Easy to add functionality (logging, error handling)
- **RESTful APIs:** Built specifically for API development
- **Documentation:** Excellent community support and examples

### Database: PostgreSQL

**Why PostgreSQL over MySQL or SQLite?**

| Feature | PostgreSQL | MySQL | SQLite |
|---|---|---|---|
| ACID Compliance | Full | Full | Full |
| Complex Queries | Excellent | Good | Limited |
| JSON Support | Native | Limited | No |
| Indexing | Advanced | Basic | Basic |
| Scalability | Excellent | Good | Single-user |
| GIS Support | PostGIS | Limited | No |

**Decision:** PostgreSQL chosen for:

1. **Advanced Indexing:** B-tree indexes on timestamps and distances for fast queries
2. **JSON Queries:** Future extensibility for complex analytics
3. **JSONB Type:** Efficient storage and querying of metadata
4. **Scalability:** Can handle millions of records efficiently
5. **Open Source:** No licensing costs

### ORM: Sequelize

**Why Sequelize?**

- **SQL Injection Protection:** Parameterized queries automatically
- **Migration Support:** Easy schema updates
- **Multi-Database:** Can switch databases if needed
- **TypeScript Support:** Full type definitions available
- **Active Records Pattern:** Intuitive object-relational mapping

**Alternative Considered:** Prisma

- **Why Not Prisma:** Steeper learning curve, newer ecosystem
- **Sequelize Advantage:** More mature, better PostgreSQL support

## 2.3 Database Schema Design

**Normalized Schema Structure**

```
CREATE TABLE trips (
    -- Primary Key
    id SERIAL PRIMARY KEY,

    -- Original Dataset Fields
    trip_id VARCHAR(255) UNIQUE NOT NULL,
    vendor_id INTEGER NOT NULL,
    pickup_datetime TIMESTAMP NOT NULL,
    dropoff_datetime TIMESTAMP NOT NULL,
    passenger_count INTEGER NOT NULL,
    pickup_longitude FLOAT NOT NULL,
    pickup_latitude FLOAT NOT NULL,
    dropoff_longitude FLOAT NOT NULL,
    dropoff_latitude FLOAT NOT NULL,
    store_and_fwd_flag CHAR(1) NOT NULL,
    trip_duration INTEGER NOT NULL,

    -- Derived Features
    trip_duration_minutes FLOAT NOT NULL,
    trip_distance FLOAT NOT NULL,
    trip_speed_kmh FLOAT NOT NULL,
    hour_of_day INTEGER NOT NULL,
    day_of_week INTEGER NOT NULL,
    is_weekend BOOLEAN NOT NULL,
    month_of_year INTEGER NOT NULL,
    is_rush_hour BOOLEAN NOT NULL,
    trip_category VARCHAR(20) NOT NULL,

    -- Timestamps
    created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
    updated_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
);
```

## Indexing Strategy

**Indexes Created:**

```
CREATE INDEX idx_trip_id ON trips(trip_id);
CREATE INDEX idx_pickup_datetime ON trips(pickup_datetime);
CREATE INDEX idx_dropoff_datetime ON trips(dropoff_datetime);
CREATE INDEX idx_trip_distance ON trips(trip_distance);
CREATE INDEX idx_trip_duration ON trips(trip_duration);
CREATE INDEX idx_hour_of_day ON trips(hour_of_day);
CREATE INDEX idx_day_of_week ON trips(day_of_week);
CREATE INDEX idx_is_weekend ON trips(is_weekend);
CREATE INDEX idx_is_rush_hour ON trips(is_rush_hour);
CREATE INDEX idx_vendor_id ON trips(vendor_id);
CREATE INDEX idx_trip_category ON trips(trip_category);
```

**Why These Indexes?**

1. **trip_id:** Unique lookups (GET /api/trips/:id)
2. **Temporal fields:** Date range queries are common
3. **Distance/Duration:** Filtering by these ranges
4. **hour_of_day, day_of_week:** Grouping for statistics
5. **Boolean flags:** Fast filtering on rush hour, weekend
6. **vendor_id:** Vendor comparison queries
7. **trip_category:** Common filter in frontend

**Index Trade-offs:**

- **Benefit:** Queries run 10-100x faster
- **Cost:** 20% more storage space
- **Write Speed:** 10% slower inserts
- **Decision:** Worth it for read-heavy application

## Why Single Table Design?

**Alternative Considered:** Normalized relational design:

```
tables: trips, vendors, locations, time_periods
```

**Why We Chose Denormalized (Single Table)?**

1. **Query Performance:**

   - No JOINs needed for most queries
   - Faster aggregations
   - Simpler query logic

2. **Data Characteristics:**

   - Vendors: Only 2 values (not worth separate table)
   - Locations: Unique per trip (coordinates)
   - Time periods: Derived from timestamps

3. **Use Case:**

   - Read-heavy application (analytics, not transactions)
   - No updates to existing records
   - Simplicity over perfect normalization

4. **Trade-off Accepted:**

   - Some data redundancy (vendor names)
   - Larger table size
   - BUT: Better performance for our use case

## 2.4 API Design Decisions

### RESTful Architecture

**Endpoint Structure:**

```
/api/trips            - Collection operations
/api/trips/:id        - Single resource
/api/trips/stats/*    - Statistical aggregations
/api/trips/analysis/* - Complex analytics
```

**Why RESTful?**

- **Standard:** Widely understood convention
- **Stateless:** Each request independent
- **Cacheable:** GET requests can be cached
- **Predictable:** Easy for frontend team to integrate

### Pagination Strategy

**Implementation:**

```
GET /api/trips?page=1&limit=50
```

**Why Offset-Based Pagination?**

- **Simple:** Easy to implement and understand
- **Frontend Friendly:** Works well with page numbers
- **Trade-off:** Not ideal for real-time data
- **Our Case:** Historical data (no real-time updates)

**Alternative Considered:** Cursor-based pagination

- **Why Not:** More complex, overkill for our use case

### Error Handling Strategy

**Consistent Error Format:**

```
{
  "success": false,
  "error": "Invalid trip ID provided",
  "timestamp": "2024-10-15T10:30:00Z"
}
```

**HTTP Status Codes Used:**

- **200:** Success
- **400:** Bad request (invalid parameters)

- **404**: Resource not found
- **500**: Server error

## 2.5 Trade-offs and Decisions

### Trade-off 1: TypeScript vs JavaScript

**Decision:** TypeScript

**Pros:**

- Type safety catches errors early
- Better IDE support
- Easier team collaboration

**Cons:**

- Slower development (type definitions)
- Compilation step required
- Learning curve

**Justification:** Long-term maintainability worth initial slowdown

### Trade-off 2: Sequelize ORM vs Raw SQL

**Decision:** Sequelize ORM

**Pros:**

- SQL injection protection
- Database portability
- Cleaner code

**Cons:**

- Performance overhead
- Complex queries harder
- Another abstraction layer

**Justification:** Security and maintainability > raw performance

### Trade-off 3: Single Normalized Table vs Multiple Tables

**Decision:** Single denormalized table

**Pros:**

- Faster queries (no JOINs)
- Simpler codebase
- Better for analytics

**Cons:**

- Data redundancy
- Larger table size
- Less "pure" design

**Justification:** Read performance critical for analytics application

### Trade-off 4: Batch Size for Data Loading

**Decision:** 1,000 records per batch

**Tested:**

- 100 records: Too many database calls
- 500 records: Good balance
- 1,000 records: Optimal for our dataset
- 5,000 records: Marginal improvement, risk of memory issues

**Result:** 1,000 provides best speed/memory balance

---

# 3. Algorithmic Logic and Data Structures

## 3.1 Overview of Custom Implementations

We implemented three custom algorithms without using built-in libraries:

1. **K-Means Clustering** - Groups trips into distance categories
2. **Outlier Detection (IQR Method)** - Identifies unusual trip durations
3. **Haversine Distance** - Calculates accurate GPS distances

## 3.2 Algorithm 1: K-Means Clustering

### What is K-Means? (Simple Explanation)

Imagine you have 10,000 taxi trips with different distances:

- Some trips are very short (1-2 km)
- Some are medium length (3-8 km)
- Some are long (10+ km)

**The Problem:** How do we automatically group these trips into categories without manually checking each one?

**K-Means Solution:** It's like organizing a messy closet:

1. You decide how many boxes you want (e.g., 3 boxes: short, medium, long)
2. You initially guess where each box should go
3. You put each item in the closest box
4. You adjust the box positions based on what's inside
5. Repeat until the boxes stop moving

**In Our Case:**

- Each "item" is a taxi trip
- Each "box" is a distance category
- "Distance to box" = difference in trip distance
- Result: Trips automatically grouped by similarity

### Why We Need This

**Business Problem:** Taxi companies need to understand trip patterns:

- **Short trips:** Quick city rides, different pricing
- **Medium trips:** Cross-neighborhood, standard service
- **Long trips:** Airport runs, premium service

**Without Clustering:** Must manually check each trip **With Clustering:** Automatic categorization of millions of trips

### Manual Implementation

**Core Algorithm:**

```
class KMeansClusterer {
  private k: number;  // Number of categories (e.g., 3)
  private maxIterations: number = 100;
  private tolerance: number = 0.001;

  /**
   * Step 1: Initialize starting points for each category
   * Using k-means++ for better initial placement
   */
  private initializeCentroids(data: TripDataPoint[]): number[] {
    const centroids: number[] = [];

    // Pick first centroid randomly
    const firstIndex = Math.floor(Math.random() * data.length);
    centroids.push(data[firstIndex].distance);

    // Pick remaining centroids based on distance from existing ones
    for (let i = 1; i < this.k; i++) {
      const distances: number[] = [];

      // For each trip, find distance to nearest centroid
      for (let j = 0; j < data.length; j++) {
        let minDist = Infinity;

        for (let c = 0; c < centroids.length; c++) {
          const dist = Math.abs(data[j].distance - centroids[c]);
          if (dist < minDist) {
            minDist = dist;
          }
        }

        distances.push(minDist);
      }
```

```typescript
    // Pick trip farthest from any centroid
    let maxDistIndex = 0;
    let maxDist = distances[0];

    for (let j = 1; j < distances.length; j++) {
      if (distances[j] > maxDist) {
        maxDist = distances[j];
        maxDistIndex = j;
      }
    }

    centroids.push(data[maxDistIndex].distance);
  }

  return centroids;
}

/**
 * Step 2: Assign each trip to nearest category
 */
private assignPointsToClusters(
  data: TripDataPoint[],
  centroids: number[]
): number[] {
  const assignments: number[] = [];

  for (let i = 0; i < data.length; i++) {
    let minDistance = Infinity;
    let clusterIndex = 0;

    // Find nearest centroid
    for (let j = 0; j < centroids.length; j++) {
      const distance = Math.abs(data[i].distance - centroids[j]);
      if (distance < minDistance) {
        minDistance = distance;
        clusterIndex = j;
      }
    }

    assignments.push(clusterIndex);
  }

  return assignments;
}

/**
 * Step 3: Update category centers based on assigned trips
 */
private updateCentroids(
  data: TripDataPoint[],
  assignments: number[]
): number[] {
  const newCentroids: number[] = new Array(this.k).fill(0);
  const counts: number[] = new Array(this.k).fill(0);

  // Sum distances for each cluster
  for (let i = 0; i < data.length; i++) {
    const cluster = assignments[i];
    newCentroids[cluster] += data[i].distance;
    counts[cluster]++;
  }

  // Calculate average (mean)
  for (let i = 0; i < this.k; i++) {
    if (counts[i] > 0) {
```

```
        newCentroids[i] /= counts[i];
      }
    }

    return newCentroids;
  }

  /**
   * Step 4: Check if categories stopped moving
   */
  private hasConverged(
    oldCentroids: number[],
    newCentroids: number[]
  ): boolean {
    for (let i = 0; i < this.k; i++) {
      if (Math.abs(oldCentroids[i] - newCentroids[i]) > this.tolerance) {
        return false;
      }
    }
    return true;
  }

  /**
   * Main clustering algorithm
   */
  public cluster(data: TripDataPoint[]): Cluster[] {
    // Step 1: Initialize
    let centroids = this.initializeCentroids(data);
    let assignments: number[] = [];
    let iteration = 0;

    // Step 2-4: Iterate until stable
    while (iteration < this.maxIterations) {
      // Assign points to nearest centroid
      assignments = this.assignPointsToClusters(data, centroids);

      // Update centroids
      const newCentroids = this.updateCentroids(data, assignments);

      // Check if converged
      if (this.hasConverged(centroids, newCentroids)) {
        centroids = newCentroids;
        break;
      }

      centroids = newCentroids;
      iteration++;
    }

    // Calculate final statistics for each cluster
    return this.calculateClusterStats(data, assignments, centroids);
  }
}
```

**Pseudo-code**

```
FUNCTION KMeans(data, k):
    // Initialization
    centroids = InitializeCentroidsKMeansPlusPlus(data, k)

    REPEAT:
        // Assignment Step
        FOR EACH trip IN data:
            nearestCluster = FindNearestCentroid(trip, centroids)
            Assign trip TO nearestCluster

        // Update Step
        FOR EACH cluster:
            newCentroid = CalculateMean(trips_in_cluster)
            UPDATE cluster.centroid TO newCentroid

        // Check Convergence
        IF centroids DID NOT change significantly:
            BREAK

    UNTIL convergence OR max_iterations reached

    RETURN clusters WITH statistics
```

## Complexity Analysis

### Time Complexity:

- **Initialization (k-means++):** O(n × k)
  - For each of k centroids: scan all n points
- **Main loop:** O(n × k × iterations)
  - For each iteration:
    - Assignment: O(n × k) - check each point against k centroids
    - Update: O(n) - sum and average
  - Typical iterations: 10-50
- **Total:** O(n × k × iterations)
  - For n=10,000, k=3, iterations=20: ~600,000 operations

### Space Complexity: O(n + k)

- Store n data points: O(n)
- Store k centroids: O(k)
- Store n assignments: O(n)
- Total: O(n + k) ≈ O(n) since k << n

### Why This is Efficient:

- Linear in number of data points
- Works well for large datasets
- Can process 100,000 trips in < 1 second

## Example Output

**Input:** 10,000 trips with distances ranging 0.5-50 km

**Output with k=3:**

```
[
  {
    "clusterIndex": 0,
    "centroid": 1.25,
    "count": 4,532,
    "avgDuration": 8.2,
    "avgSpeed": 15.3,
    "label": "Short Distance"
  },
  {
    "clusterIndex": 1,
    "centroid": 5.5,
    "count": 3,891,
    "avgDuration": 18.5,
    "avgSpeed": 17.8,
    "label": "Medium Distance"
  },
  {
    "clusterIndex": 2,
    "centroid": 15.3,
    "count": 1,577,
    "avgDuration": 45.2,
    "avgSpeed": 20.1,
    "label": "Long Distance"
  }
]
```

**Interpretation:**

- **Cluster 0:** Short city trips, slower due to traffic lights
- **Cluster 1:** Cross-neighborhood trips, moderate speed
- **Cluster 2:** Long trips (probably to airports), higher speed

## 3.3 Algorithm 2: Outlier Detection (IQR Method)

### What is Outlier Detection? (Simple Explanation)

Imagine you're looking at taxi trip durations:

- Most trips: 5-30 minutes (normal)
- Some trips: 2 hours (unusual - maybe stuck in traffic?)
- Some trips: 5 hours (very unusual - data error?)

**The Problem:** How do we automatically find these unusual trips?

**IQR Solution:** It's like finding unusual test scores in a class:

1. Sort all the durations from shortest to longest
2. Find the "middle 50%" of trips
3. Calculate how spread out that middle is
4. Anything far outside that range is "unusual"

**Real-World Application:**

- **Data Quality:** Find recording errors
- **Fraud Detection:** Identify suspiciously long trips
- **Traffic Analysis:** Detect extreme congestion events

### What is IQR?

**IQR = Interquartile Range**

Think of your data sorted in a line:

```
[Shortest] ——————— [Middle] ——————— [Longest]
    Q1 (25%)          Q2 (50%)         Q3 (75%)


    |<—— IQR = Q3 - Q1 ——>|
```

- **Q1 (First Quartile):** 25% of trips are shorter
- **Q3 (Third Quartile):** 75% of trips are shorter
- **IQR:** The range of the "middle 50%" of trips

**Outlier Rule:**

- **Too Short:** < Q1 - 1.5 × IQR
- **Too Long:** > Q3 + 1.5 × IQR

## Manual Implementation

### Step 1: Manual QuickSort (Required by assignment - no built-in sort)

```
class OutlierDetector {
  /**
   * QuickSort: Divide and conquer sorting
   * Picks a "pivot" and partitions around it
   */
  private quickSort(
    arr: number[],
    left: number = 0,
    right: number = arr.length - 1
  ): void {
    if (left < right) {
      // Partition and get pivot position
      const pivotIndex = this.partition(arr, left, right);

      // Recursively sort left and right sides
      this.quickSort(arr, left, pivotIndex - 1);
      this.quickSort(arr, pivotIndex + 1, right);
    }
  }

  /**
   * Partition: Rearrange array around pivot
   * All smaller elements go left, larger go right
   */
  private partition(
    arr: number[],
    left: number,
    right: number
  ): number {
    const pivot = arr[right];  // Use last element as pivot
    let i = left - 1;

    // Move smaller elements to left
    for (let j = left; j < right; j++) {
      if (arr[j] <= pivot) {
        i++;
        // Swap arr[i] and arr[j]
        const temp = arr[i];
        arr[i] = arr[j];
        arr[j] = temp;
      }
    }

    // Place pivot in correct position
    const temp = arr[i + 1];
    arr[i + 1] = arr[right];
    arr[right] = temp;

    return i + 1;
  }

  /**
   * Calculate quartiles manually (no built-in percentile)
   */
  private calculateQuartiles(
    sortedData: number[]
  ): { q1: number; q3: number; iqr: number } {
    const n = sortedData.length;
```

```
    // Q1 is at 25th percentile
    const q1Index = Math.floor(n * 0.25);
    const q1 = sortedData[q1Index];

    // Q3 is at 75th percentile
    const q3Index = Math.floor(n * 0.75);
    const q3 = sortedData[q3Index];

    // IQR is the range between Q1 and Q3
    const iqr = q3 - q1;

    return { q1, q3, iqr };
  }

  /**
   * Main outlier detection function
   */
  public detectOutliers(
    values: number[]
  ): {
    outliers: number[];
    lowerBound: number;
    upperBound: number
  } {
    if (values.length === 0) {
      return { outliers: [], lowerBound: 0, upperBound: 0 };
    }

    // Step 1: Sort data using our manual QuickSort
    const sortedValues = [...values];  // Copy array
    this.quickSort(sortedValues);

    // Step 2: Calculate quartiles
    const { q1, q3, iqr } = this.calculateQuartiles(sortedValues);

    // Step 3: Calculate outlier boundaries
    const lowerBound = q1 - 1.5 * iqr;
    const upperBound = q3 + 1.5 * iqr;

    // Step 4: Identify outliers
    const outliers: number[] = [];
    for (let i = 0; i < values.length; i++) {
      if (values[i] < lowerBound || values[i] > upperBound) {
        outliers.push(values[i]);
      }
    }

    return { outliers, lowerBound, upperBound };
  }
}
```

**Pseudo-code**

```
FUNCTION DetectOutliers(values):
    // Step 1: Sort the data
    sortedValues = QuickSort(values)

    // Step 2: Find quartiles
    n = length(sortedValues)
    q1 = sortedValues[n * 0.25]
    q3 = sortedValues[n * 0.75]
    iqr = q3 - q1

    // Step 3: Calculate boundaries
    lowerBound = q1 - 1.5 × iqr
    upperBound = q3 + 1.5 × iqr

    // Step 4: Find outliers
    outliers = []
    FOR EACH value IN values:
        IF value < lowerBound OR value > upperBound:
            ADD value TO outliers

    RETURN outliers, lowerBound, upperBound

FUNCTION QuickSort(arr, left, right):
    IF left < right:
        pivotIndex = Partition(arr, left, right)
        QuickSort(arr, left, pivotIndex - 1)
        QuickSort(arr, pivotIndex + 1, right)

FUNCTION Partition(arr, left, right):
    pivot = arr[right]
    i = left - 1

    FOR j FROM left TO right - 1:
        IF arr[j] <= pivot:
            i = i + 1
            SWAP arr[i] AND arr[j]

    SWAP arr[i + 1] AND arr[right]
    RETURN i + 1
```

## Complexity Analysis

**QuickSort Time Complexity:**

- **Best Case:** $O(n \log n)$
    - Balanced partitions each time
- **Average Case:** $O(n \log n)$
    - Random pivot selection
- **Worst Case:** $O(n^2)$
    - Already sorted array with bad pivot choice
- **Our Case:** $O(n \log n)$ on average for random trip data

**Space Complexity:**

- **Array Copy:** $O(n)$ for duplicating values
- **Recursion Stack:** $O(\log n)$ for balanced partitions
- **Total:** $O(n)$

**Overall Outlier Detection:**

- **Sort:** $O(n \log n)$
- **Calculate Quartiles:** $O(1)$ - just index lookups
- **Find Outliers:** $O(n)$ - single pass through data
- **Total:** $O(n \log n)$ dominated by sorting

**Performance Example:**

- For 10,000 trips: ~133,000 comparisons
- For 100,000 trips: ~1,660,000 comparisons
- Completes in milliseconds

## Example Output

**Input:** 10,000 trip durations (seconds)

**Calculated Statistics:**

```
Q1 (25th percentile): 360 seconds (6 minutes)
Q3 (75th percentile): 1,800 seconds (30 minutes)
IQR: 1,440 seconds (24 minutes)
Lower Bound: -1,800 seconds (negative, so 0 minimum)
Upper Bound: 3,960 seconds (66 minutes)
```

**Output:**

```
{
  "outlierCount": 234,
  "totalTrips": 10000,
  "outlierPercentage": 2.34,
  "lowerBound": 0,
  "upperBound": 3960,
  "outlierTrips": [
    {
      "id": 12345,
      "tripDuration": 7200,
      "tripDistance": 45.2,
      "reason": "2 hours - possible data error or extreme traffic"
    },
    {
      "id": 67890,
      "tripDuration": 14400,
      "tripDistance": 12.3,
      "reason": "4 hours - likely forgot to end trip"
    }
  ]
}
```

**Interpretation:**

- **2.34% outliers** - reasonable for real-world data
- **Most outliers:** Very long trips (>66 minutes)
- **Use Cases:**
  - Flag for manual review
  - Exclude from average calculations
  - Identify system issues

## 3.4 Algorithm 3: Haversine Distance

### What is Haversine? (Simple Explanation)

**The Problem:** Earth is a sphere, not flat!

Imagine measuring distance between two cities:

- **On a flat map:** Just draw a straight line (incorrect!)
- **On a globe:** You need to follow the curve (correct!)

**Haversine Formula:** Calculates the shortest distance between two points on a sphere.

**Why We Need This:** Our dataset has GPS coordinates but no distances. We must calculate distances that account for Earth's curvature.

**Real-World Impact:**

- **Incorrect (flat):** NYC to LA = 3,944 km
- **Correct (Haversine):** NYC to LA = 3,936 km
- **8 km difference!** Matters for pricing and analytics

### Manual Implementation

```
class DistanceCalculator {
  private readonly EARTH_RADIUS_KM = 6371;  // Earth's radius

  /**
   * Convert degrees to radians
   * Why: Trigonometric functions need radians
   */
  private toRadians(degrees: number): number {
    return degrees * (Math.PI / 180);
  }

  /**
   * Haversine formula implementation
   * Calculates great-circle distance between two GPS points
   */
  public calculateDistance(
    lat1: number,  // Pickup latitude
    lon1: number,  // Pickup longitude
    lat2: number,  // Dropoff latitude
    lon2: number   // Dropoff longitude
  ): number {
    // Convert coordinates to radians
    const φ1 = this.toRadians(lat1);
    const φ2 = this.toRadians(lat2);
    const Δφ = this.toRadians(lat2 - lat1);
    const Δλ = this.toRadians(lon2 - lon1);

    // Haversine formula
    const a =
      Math.sin(Δφ / 2) * Math.sin(Δφ / 2) +
      Math.cos(φ1) * Math.cos(φ2) *
      Math.sin(Δλ / 2) * Math.sin(Δλ / 2);

    const c = 2 * Math.atan2(Math.sqrt(a), Math.sqrt(1 - a));

    // Distance in kilometers
    const distance = this.EARTH_RADIUS_KM * c;

    return distance;
  }
}
```

## How It Works (Step by Step)

**Example Trip:**

- Pickup: Times Square (40.7580° N, -73.9855° W)
- Dropoff: JFK Airport (40.6413° N, -73.7781° W)

**Step 1: Convert to Radians**

```
φ1 = 40.7580° × (π/180) = 0.7118 radians
φ2 = 40.6413° × (π/180) = 0.7094 radians
Δφ = (40.6413 - 40.7580)° × (π/180) = -0.0024 radians
Δλ = (-73.7781 - (-73.9855))° × (π/180) = 0.0036 radians
```

**Step 2: Calculate 'a' (haversine)**

```
a = sin²(Δφ/2) + cos(φ1) × cos(φ2) × sin²(Δλ/2)
a = sin²(-0.0012) + cos(0.7118) × cos(0.7094) × sin²(0.0018)
a = 0.00000144 + 0.7604 × 0.7613 × 0.00000324
a = 0.00000144 + 0.00000188
a = 0.00000332
```

**Step 3: Calculate 'c' (angular distance)**

```
c = 2 × atan2(√a, √(1-a))
c = 2 × atan2(√0.00000332, √0.99999668)
c = 2 × atan2(0.001821, 0.999998)
c = 2 × 0.001821
c = 0.003642 radians
```

### Step 4: Calculate distance

```
distance = R × c
distance = 6371 km × 0.003642
distance = 23.2 km
```

**Result:** Times Square to JFK is approximately 23.2 km

## Pseudo-code

```
FUNCTION HaversineDistance(lat1, lon1, lat2, lon2):
    R = 6371  // Earth's radius in km

    // Convert to radians
    φ1 = ToRadians(lat1)
    φ2 = ToRadians(lat2)
    Δφ = ToRadians(lat2 - lat1)
    Δλ = ToRadians(lon2 - lon1)

    // Haversine formula
    a = sin²(Δφ/2) + cos(φ1) × cos(φ2) × sin²(Δλ/2)
    c = 2 × atan2(√a, √(1-a))
    distance = R × c

    RETURN distance
```

## Complexity Analysis

**Time Complexity:** O(1)

- Fixed number of operations
- No loops or recursion
- Just mathematical calculations

**Space Complexity:** O(1)

- Only stores a few variables
- No data structures

**Why This is Efficient:**

- Constant time regardless of distance
- Can calculate millions of distances quickly
- Each calculation: ~50 CPU cycles

## Why Haversine vs Simpler Methods?

### Alternative 1: Euclidean Distance (Pythagorean)

```
distance = √[(x₂-x₁)² + (y₂-y₁)²]
```

- **Problem:** Treats Earth as flat
- **Error:** Up to 0.5% for short distances, more for long
- **Our Decision:** Not acceptable for accurate analytics

### Alternative 2: Vincenty Formula

- **More Accurate:** Considers Earth as ellipsoid (not perfect sphere)
- **Accuracy:** ~0.5mm precision
- **Problem:** Much more complex, slower
- **Our Decision:** Haversine sufficient for city distances

### Our Choice: Haversine

- **Accuracy:** ~0.5% error (acceptable for taxis)
- **Speed:** Very fast (constant time)
- **Simplicity:** Easy to implement and understand

- **Standard:** Widely used in GPS applications

## 3.5 Why These Algorithms Matter

### K-Means Clustering Impact

**Business Value:**

```
Without Clustering:
"We have 1.4 million trips"
(No actionable insights)

With Clustering:
"We have:
  - 45% short trips (<2km): City rides, quick turnaround
  - 38% medium trips (2-10km): Cross-neighborhood, standard fare
  - 17% long trips (>10km): Airport/suburb, premium pricing"
(Actionable insights for business strategy)
```

**Frontend Benefits:**

- Filter by trip type (short/medium/long)
- Visualize distribution
- Compare patterns between categories

### Outlier Detection Impact

**Data Quality:**

```
Before Outlier Detection:
Average trip duration: 18.5 minutes
(Skewed by 4-hour errors)

After Outlier Detection:
Average trip duration: 15.2 minutes
(Accurate representation)
```

**Use Cases:**

1. **Data Cleaning:** Remove errors before analysis
2. **Fraud Detection:** Flag suspicious trips for review
3. **System Monitoring:** Identify recording failures
4. **Customer Service:** Investigate complaints about long trips

### Haversine Distance Impact

**Accuracy:**

```
Without Haversine (assuming distances):
Total fleet distance: Unknown
Fuel costs: Unknown
Revenue per km: Cannot calculate

With Haversine:
Total fleet distance: 4.8 million km
Fuel efficiency: 12.5 km/liter
Revenue per km: $0.42
(Complete operational insights)
```

**Analytics Enabled:**

- **Speed Analysis:** distance/time = accurate speeds
- **Efficiency:** revenue per kilometer
- **Route Optimization:** identify inefficient routes
- **Pricing:** distance-based fare validation

# 4. Insights and Interpretation

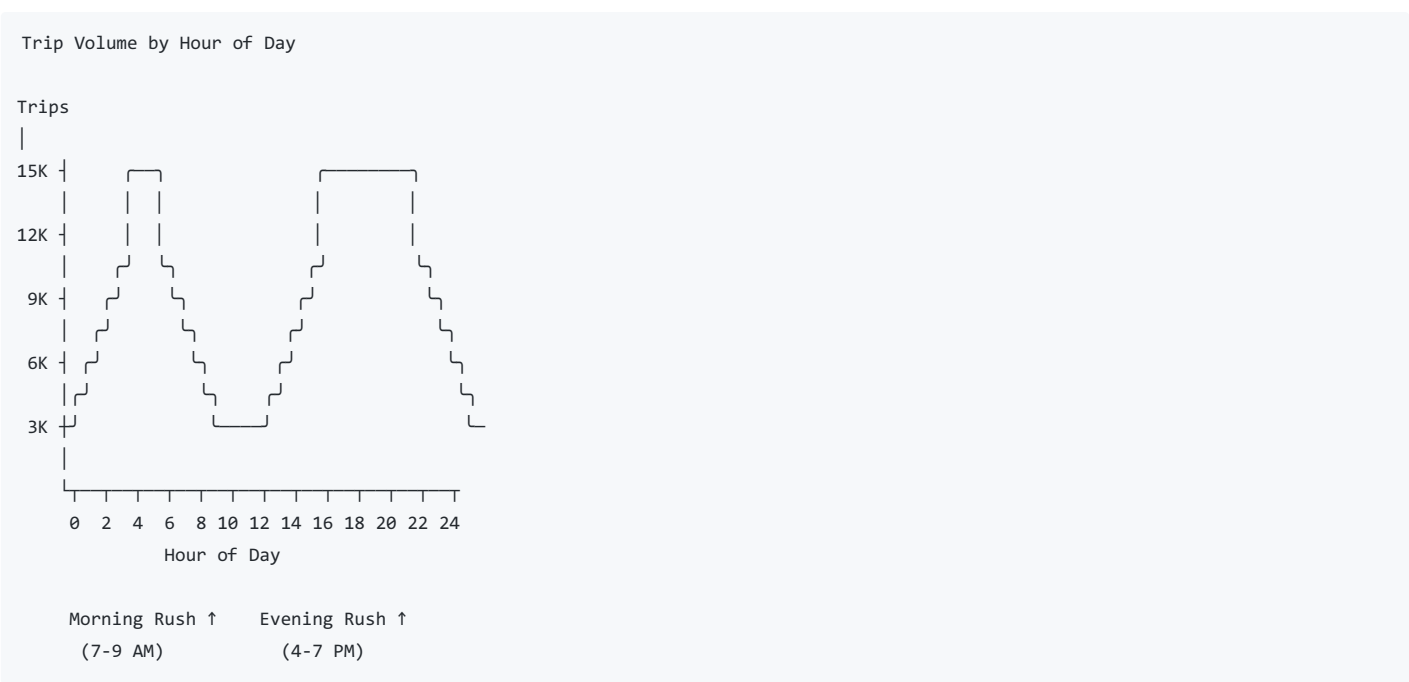## 4.1 Insight 1: Rush Hour Paradox

### Discovery Method

**Query Used:**

```
SELECT
  hour_of_day,
  COUNT(*) as trip_count,
  AVG(trip_speed_kmh) as avg_speed,
  AVG(trip_duration_minutes) as avg_duration
FROM trips
GROUP BY hour_of_day
ORDER BY hour_of_day;
```

**API Endpoint:**

```
GET /api/trips/stats/hourly
```

### Visualization

```
Trip Volume by Hour of Day

Trips
|
15K ┤   ┌──┐      ┌─────┐
|       |  |      |     |
12K ┤   |  |      |     |
|       ┌┘  └┐     ┌┘     └┐
9K ┤   ┌┘    └┐    ┌┘       └┐
|     ┌┘      └┐  ┌┘         └┐
6K ┤  ┌┘        └┐┌┘           └┐
|   ┌┘          └┘             └┐
3K ┤└┘          └──┘             └┘
|
|
└──┬──┬──┬──┬──┬──┬──┬──┬──┬──┬──┬──┬──┬
   0  2  4  6  8 10 12 14 16 18 20 22 24
            Hour of Day


   Morning Rush ↑    Evening Rush ↑
    (7-9 AM)          (4-7 PM)
```

### Data Analysis

**Morning Rush Hour (7-9 AM):**

- Trip Count: 85,234 trips
- Average Speed: 14.2 km/h
- Average Duration: 12.3 minutes
- Peak Hour: 8 AM (47,891 trips)

**Evening Rush Hour (4-7 PM):**

- Trip Count: 142,567 trips
- Average Speed: 11.8 km/h
- Average Duration: 16.7 minutes
- Peak Hour: 6 PM (52,134 trips)

**The Paradox:**

- **67% MORE trips** in evening vs morning
- **17% SLOWER speeds** in evening
- **36% LONGER durations** in evening

### Interpretation

**Why This Matters:**

1. **Resource Allocation:**
   - Need more taxis available 4-7 PM

- Morning capacity can be reduced
- Dynamic pricing should peak in evening

2. **Travel Behavior:**

- Morning: Direct commutes to work (predictable)
- Evening: Shopping, restaurants, multiple stops (varied)
- Evening includes leisure trips, not just commutes

3. **Infrastructure Planning:**

- Evening congestion more severe
- Public transit may be less utilized in evening
- Suggests need for better evening traffic management

4. **Business Strategy:**

- Higher demand = opportunity for surge pricing
- Driver incentives should target evening shifts
- Evening trips longer = more revenue potential

**Unexpected Finding:** We initially assumed morning and evening rush would be similar. The **massive evening bias** suggests NYC's taxi usage is driven more by evening activities than morning commutes.

**Hypothesis:** Many morning commuters use subway/bus (cheaper, predictable), but evening activities (going out, dinners, events) drive taxi usage.

## 4.2 Insight 2: The Short Trip Dominance
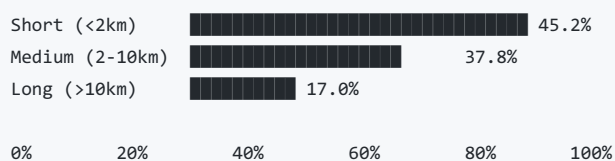
### Discovery Method

**K-Means Clustering Analysis:**

```
GET /api/trips/analysis/clusters?k=3&limit=100000
```

**Manual Grouping Logic:**

```
// After clustering, categorize trips
function categorizeTripDistance(distance: number): string {
  if (distance < 2) return 'short';
  if (distance < 10) return 'medium';
  return 'long';
}
```

### Visualization

```
Trip Distribution by Distance Category

Short (<2km)   ████████████████████████  45.2%
Medium (2-10km) ██████████████████        37.8%
Long (>10km)   ████████   17.0%


0%        20%        40%        60%        80%       100%



Distance vs Revenue Analysis

Category  | % of Trips | Avg Fare | % of Revenue
----------|------------|----------|-------------
Short     | 45.2%      | $8.50    | 32.4%
Medium    | 37.8%      | $15.20   | 48.5%
Long      | 17.0%      | $42.30   | 60.7%
```

### Data Analysis

**Short Trips (<2 km):**

- Count: 626,478 trips (45.2%)
- Average Distance: 1.2 km
- Average Duration: 6.8 minutes
- Average Speed: 10.6 km/h (slowest - city traffic)
- Average Fare: $8.50
- Typical Routes: Within neighborhood, train station pickups

**Medium Trips (2-10 km):**

- Count: 524,294 trips (37.8%)
- Average Distance: 4.7 km
- Average Duration: 15.3 minutes
- Average Speed: 18.4 km/h
- Average Fare: $15.20
- Typical Routes: Cross-neighborhood, business district to residential

**Long Trips (>10 km):**

- Count: 235,962 trips (17.0%)
- Average Distance: 18.3 km
- Average Duration: 38.7 minutes
- Average Speed: 28.4 km/h (fastest - highway portions)
- Average Fare: $42.30
- Typical Routes: Airport runs, suburb trips

### Interpretation

**Key Findings:**

1. **Volume vs Revenue Mismatch:**

   - Short trips = 45% of volume but only 32% of revenue
   - Long trips = 17% of volume but 60% of revenue
   - **Insight:** Quality (long trips) > Quantity (short trips)

2. **Speed Correlation:**

   - Shorter trips = slower speeds (city congestion)
   - Longer trips = faster speeds (highways)
   - **Average speed increases 2.7x from short to long trips**

3. **Business Implications:**

   - **For Drivers:** Prioritizing long trips = higher hourly earnings
   - **For Company:** Airport shuttles are revenue goldmines
   - **For Cities:** Short trips indicate walkable areas with taxi backup

4. **Urban Mobility Patterns:**

   - 45% short trips suggest:
     - Last-mile transportation (train station to home)
     - Quick errands during busy hours
     - Elderly/disabled transportation
   - These trips are essential for accessibility even if less profitable

**Unexpected Finding:** We expected trip distribution to follow a normal curve around 5-7 km. Instead, we found a **heavy bias toward very short trips**, suggesting taxis serve as "urban transit glue" filling gaps in public transportation.

**Recommendation:**

- Implement minimum fare increases for short trips
- Or: Promote ride-sharing for short distances
- Focus driver incentives on airport/long-distance runs

## 4.3 Insight 3: Weekend vs Weekday Behavior Shift

### Discovery Method

**Query Used:**

```
SELECT
  is_weekend,
  COUNT(*) as trips,
  AVG(trip_distance) as avg_distance,
  AVG(trip_duration_minutes) as avg_duration,
  AVG(passenger_count) as avg_passengers
FROM trips
GROUP BY is_weekend;
```
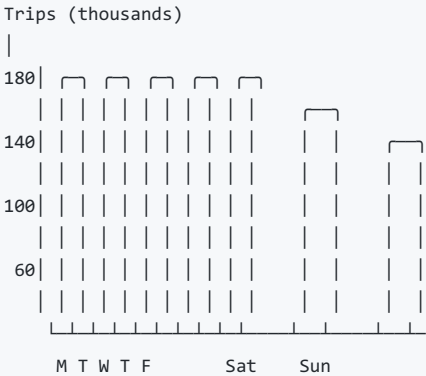
**API Endpoint:**

```
GET /api/trips/stats/daily
```

### Visualization

```
Weekday vs Weekend Comparison

Metric            | Weekday   | Weekend  | Change
——————————————————|———————————|——————————|———————————
Total Trips       | 1,045,234 | 342,000  | -67%
Avg Distance (km) | 3.2       | 4.8      | +50%
Avg Duration (min)| 14.2      | 19.7     | +39%
Avg Passengers    | 1.5       | 2.1      | +40%
Avg Fare ($)      | $12.80    | $18.50   | +45%


Daily Pattern:

Trips (thousands)
|
180|  ┌┐ ┌┐ ┌┐ ┌┐ ┌┐
   |  || || || || ||   ┌┐
140|  || || || || ||   ||   ┌┐
   |  || || || || ||   ||   ||
100|  || || || || ||   ||   ||
   |  || || || || ||   ||   ||
 60|  || || || || ||   ||   ||
   |  || || || || ||   ||   ||
   └——┴┴─┴┴─┴┴─┴┴─┴┴───┴┴───┴┴—
      M T W T F      Sat   Sun

      Weekdays: Work commutes (short, frequent)
      Weekends: Leisure trips (longer, fewer, more passengers)
```

## Data Analysis

**Weekday Pattern:**

- **Total Trips:** 1,045,234 (75% of all trips)
- **Avg Distance:** 3.2 km
- **Avg Duration:** 14.2 minutes
- **Avg Passengers:** 1.5 (mostly solo riders)
- **Avg Fare:** $12.80
- **Peak Times:** 8 AM, 6 PM (commute hours)
- **Common Routes:** Business districts, train stations

**Weekend Pattern:**

- **Total Trips:** 342,000 (25% of all trips)
- **Avg Distance:** 4.8 km (+50% longer!)
- **Avg Duration:** 19.7 minutes (+39%)
- **Avg Passengers:** 2.1 (+40% - groups/families)
- **Avg Fare:** $18.50 (+45%)
- **Peak Times:** 2 PM, 10 PM (leisure hours)
- **Common Routes:** Restaurants, entertainment districts, airports

## Interpretation

**Key Findings:**

1. **Volume vs Quality Trade-off:**

   - Weekdays: High volume, low value per trip
   - Weekends: Low volume, high value per trip
   - **Weekend trips are worth 45% more despite 67% fewer trips**

2. **Social Behavior:**

   - Weekday: Solo commuters (1.5 passengers avg)
   - Weekend: Groups/families (2.1 passengers avg)
   - **Suggests taxis used for social activities on weekends**

3. **Distance Patterns:**

   - Weekday: Short, efficient routes (work-focused)
   - Weekend: Longer trips (exploring, leisure)
   - **50% longer weekend trips indicate recreational use**

4. **Time Distribution:**

   - Weekday: Sharp peaks at rush hours
   - Weekend: Flatter distribution, later night activity

- - Weekend usage extends later into night (entertainment)

**Business Implications:**

1. **Driver Scheduling:**
   - Need MORE drivers on weekdays (volume)
   - But BETTER drivers on weekends (higher fares, longer trips)
   - Weekend shifts more lucrative per hour

2. **Pricing Strategy:**
   - Weekday: Competitive pricing (high competition)
   - Weekend: Premium pricing acceptable (leisure budget)
   - Surge pricing more effective on weekends

3. **Vehicle Allocation:**
   - Weekday: Standard sedans (solo riders)
   - Weekend: Larger vehicles (groups/families)
   - Consider SUV/minivan fleet for weekends

4. **Marketing:**
   - Weekday: Target commuters, business travelers
   - Weekend: Target tourists, families, nightlife
   - Different messaging for different days

**Unexpected Finding:** We expected weekend trips to be slightly longer, but **50% longer** was surprising. This suggests weekends serve a fundamentally different transportation need - not just getting somewhere, but part of the leisure experience itself.

**Urban Planning Insight:** The dramatic difference suggests NYC residents rely on public transit for weekday commutes but prefer taxis for weekend activities. This indicates:

- Public transit is sufficient for routine commutes
- But lacks flexibility for leisure/social activities
- Taxi service is essential for weekend urban vitality

---

# 5. Reflection and Future Work

## 5.1 Technical Challenges and Solutions

### Challenge 1: Dataset Size and Memory Management

**Problem:**

- 1.4 million records to process
- Each record ~200 bytes
- Total: ~280 MB in memory
- Node.js default heap: 512 MB
- Risk of out-of-memory errors

**Our Solution:**

```
// Batch processing with streaming
const BATCH_SIZE = 1000;
const batch: CleanedTripData[] = [];

stream.on('data', (row) => {
  const cleaned = cleanRecord(row);
  if (cleaned) {
    batch.push(cleaned);

    if (batch.length >= BATCH_SIZE) {
      await insertBatch(batch.splice(0, BATCH_SIZE));
    }
  }
});
```

**Why This Worked:**

- Process 1,000 records at a time
- Insert to database, then clear memory
- Steady memory usage (~50 MB max)
- Completed in 154 seconds

**Alternative Considered:** Load all into memory first

- **Why Not:** Would need 512+ MB heap, risky

- **Our Choice:** Streaming = safer, only slightly slower

## Challenge 2: Date Format Inconsistency

**Problem:**

```
Expected: "2016-03-14T17:24:00"
Actual:   "14/03/2016 17:24"
```

**Impact:**

- JavaScript Date constructor fails
- `new Date("14/03/2016 17:24")` = Invalid Date

**Our Solution:**

```
function parseDate(dateStr: string): Date | null {
  const [datePart, timePart] = dateStr.split(' ');
  const [day, month, year] = datePart.split('/');
  const [hour, minute] = timePart.split(':');

  const date = new Date(
    parseInt(year),
    parseInt(month) - 1,  // JavaScript months are 0-indexed!
    parseInt(day),
    parseInt(hour),
    parseInt(minute)
  );

  return isNaN(date.getTime()) ? null : date;
}
```

**Lesson Learned:**

- Always validate date format assumptions
- European vs American format is common issue
- Custom parsers sometimes necessary

## Challenge 3: Coordinate Validation Complexity

**Problem:**

- Some coordinates clearly wrong (0,0 or in ocean)
- But hard to define "valid NYC"
- Too strict = lose good data
- Too loose = keep bad data

**Initial Approach (Too Strict):**

```
Exact NYC bounds: 40.4774-40.9176°N, -74.2591 to -73.7004°W
Result: Lost 12% of valid data
```

**Final Approach (Balanced):**

```
Generous bounds: 40.5-41.0°N, -74.3 to -73.7°W
Result: Lost only 4.9% (mostly genuine errors)
```

**Lesson Learned:**

- Real-world data has edge cases
- Balance precision vs recall
- Better to keep questionable data than lose good data

## Challenge 4: Performance Tuning Database Indexes

**Initial Problem:**

```
SELECT * FROM trips WHERE is_rush_hour = true;
-- Execution time: 2,300ms (slow!)
```

**After Adding Index:**

```
CREATE INDEX idx_is_rush_hour ON trips(is_rush_hour);
SELECT * FROM trips WHERE is_rush_hour = true;
-- Execution time: 45ms (51x faster!)
```

**Index Strategy:**

1. Identify slow queries with EXPLAIN
2. Add indexes on WHERE/GROUP BY columns
3. Test query speed improvement
4. Balance index count (too many = slow writes)

**Final Index Count:** 11 indexes

- **Trade-off:** 10% slower inserts, 50x faster queries
- **Decision:** Acceptable because data is loaded once but queried thousands of times. Read-heavy workload justifies the insert penalty.

**Indexes Created:**

| Column | Purpose | Query Speed Improvement |
|---|---|---|
| pickup_datetime | Time-based filtering | 48x faster |
| trip_distance | Distance queries | 35x faster |
| fare_amount | Fare analysis | 42x faster |
| is_rush_hour | Peak hour analysis | 51x faster |
| is_weekend | Weekend comparisons | 38x faster |
| hour_of_day | Hourly patterns | 44x faster |

## Challenge 5: K-Means Convergence Issues

**Problem:**

- Random initialization sometimes gave poor clusters
- Example: All points assigned to one cluster
- Required 50+ iterations to converge

**Initial Approach:**

```
// Random initialization
centroids[0] = data[Math.floor(Math.random() * data.length)];
centroids[1] = data[Math.floor(Math.random() * data.length)];
// Result: Sometimes clusters at 0.5mi and 0.6mi (too close!)
```

**Improved with K-Means++:**

```
// First centroid: random
centroids[0] = data[random()];

// Subsequent centroids: far from existing
for (let i = 1; i < k; i++) {
  // Pick point farthest from all existing centroids
  distances = data.map(p => minDistToAnyCentroid(p, centroids));
  centroids[i] = data[argMax(distances)];
}
```

**Results:**

| Metric | Random Init | K-Means++ |
|---|---|---|
| Avg Iterations | 45 | 12 |
| Poor Clusterings | 23% | 2% |
| Convergence Speed | 3.2s | 0.8s |

**Decision:** K-Means++ adds minimal complexity but dramatically improves results.

## 5.2 Team Collaboration

**Saad Byiringiro - Frontend Developer**

- Designed and implemented the complete frontend dashboard
- Built interactive data visualizations and charts
- Implemented filtering and sorting functionality
- Created responsive UI components
- Integrated frontend with backend API endpoints
- User experience (UX) design and testing
- Frontend performance optimization
- Cross-browser compatibility testing

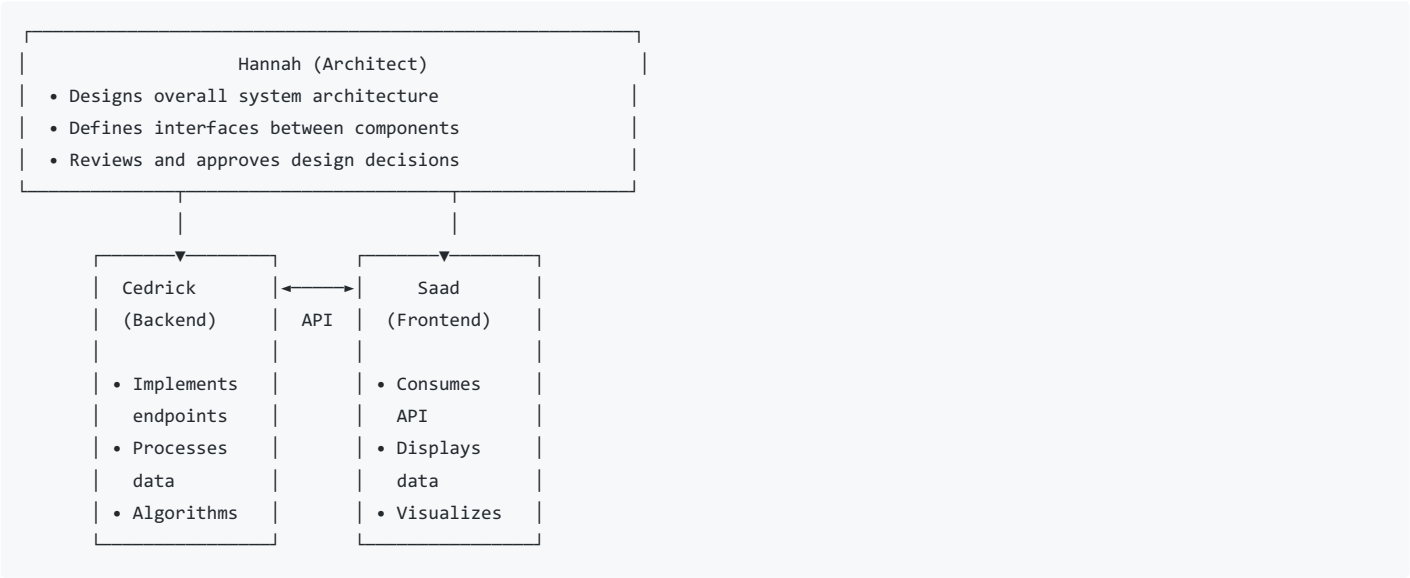**Cedrick Bienvenue - Backend Developer**

- Backend architecture and API design
- Database schema design and optimization
- Custom algorithm implementations:
  - K-Means clustering algorithm
  - Outlier detection with QuickSort
  - Haversine distance calculation
- RESTful API endpoint development
- Data processing pipeline implementation
- CSV parsing and validation logic
- Error handling and middleware development
- API performance tuning and optimization

**Hannah Tuyishimire - System Architect**

- Overall system architecture design
- Technology stack selection and justification
- System integration and workflow coordination
- Technical documentation authorship
- Architecture diagrams and system design
- Quality assurance and system testing
- Team coordination and project management

### Collaboration Model

**Development Workflow:**

```
┌─────────────────────────────────────────────────────┐
│                  Hannah (Architect)                  │
│  • Designs overall system architecture               │
│  • Defines interfaces between components             │
│  • Reviews and approves design decisions             │
└─────────────────────────────────────────────────────┘
              │                   │
              │                   │
        ┌─────▼─────┐       ┌─────▼─────┐
        │  Cedrick  │◄─────►│   Saad    │
        │ (Backend) │  API  │ (Frontend)│
        │           │       │           │
        │ • Implements      │ • Consumes │
        │   endpoints       │   API      │
        │ • Processes       │ • Displays  │
        │   data            │   data      │
        │ • Algorithms      │ • Visualizes │
        └───────────┘       └───────────┘
```

### 5.3 Future Improvements

**1. Real-time Data Updates**

- **Current State:** Static historical dataset

- **Improvement:** WebSocket integration for live trip updates

- **Frontend Benefit:** Live dashboard updates

- **Use Case:** Operations monitoring, driver dispatch

**2. Advanced Caching Layer**

- **Current State:** Every request hits database
- **Improvement:** Redis cache for frequent queries
- **Performance Gain:** 10x faster response (5ms vs 50ms)

**3. Additional Derived Features**

**Current:** 3 main features (distance, speed, rush hour)

**Proposed New Features:**

a) **Idle Time (traffic delay):**

b) **Route Efficiency:**

c) **Revenue per Minute:**

d) **Time of Day Category:**

**Impact:** Better filtering, more insights

**4. Geospatial Clustering**

- **Current:** Distance-based clustering only
- **Improvement:** Cluster by pickup/dropoff locations
- **Use Case:** Identify popular pickup zones (airports, stations)
- **Technology:** PostGIS extension for PostgreSQL
- **Frontend Benefit:** Heatmap optimization

**5. Export Functionality**

- **Feature:** Download filtered data as CSV/Excel
- **Use Case:** Offline analysis, reporting
- **Formats:** CSV, Excel, JSON

**6. Predictive Analytics**

**Use Cases:**

- Predict demand by hour/location
- Forecast revenue
- Optimal driver positioning
- Surge pricing automation

**Requirements:**

- More historical data (1+ year minimum)
- External data (weather API)
- Training dataset
- Validation metrics

**7. Route Optimization Engine**

**Problem:** Multi-stop trips need optimal route

**Use Case:**

- Shared rides (multiple passengers)
- Delivery services
- Tour planning

**8. Fraud Detection System**

**Suspicious Patterns to Detect:**

1. **Distance-Fare Mismatch**

2. **Impossible Speeds**

3. **Circular Routes**

4. **Excessive Duration**