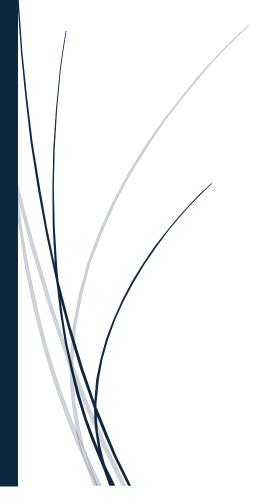01.04.2025

# 6.1 Data Research Project
Logistics Route Optimization in Estonia

by Denis Kleptsov
CAREER FOUNDRY DATA ANALYTICS COURSE

Contents

# Introduction

This project aims to optimize delivery routes for a textile logistics operation in Estonia using clustering techniques and route analysis. By grouping customers geographically and refining delivery routes, the company seeks to reduce travel distance and time while maintaining service quality. Two new datasets – one detailing individual customer deliveries and another summarizing route distances and times – provide the basis for analysis. The focus is on clustering delivery locations, optimizing route sequences, and analyzing correlations between delivery volumes and route performance. All analyses and plans align to improve operational efficiency within the data scope.

# 1  Project Alignment

The project **objectives** are:

- **Cluster Customers**: Group customers by geographic proximity to reduce travel distance by at least 15%. Develop efficient delivery routes within each cluster, minimizing distance and time while adhering to vehicle load limits, service times, working hours, and depot service times, followed by determining the number of new routes and their stops.
- **Optimize Routes**: Develop efficient delivery routes within each cluster, minimizing distance and time while adhering to vehicle load limits, service times, working hours, and depot service times, followed by determining the number of new routes and their stops.
- **Rationale**: Incorporates constraints explicitly, ensuring routes are practical and measurable against KPIs.

**Scope**

- **Geographical Focus**: Estonia, based on the provided datasets.
- **Timeframe**: The available data covers one week (March 17–23, 2025). This snapshot will be used for clustering and route optimization analysis, recognizing that it represents a short-term operational window.
- **Exclusions**: Due to data unavailability, factors such as real-time traffic conditions, weather impacts, and specific vehicle constraints (e.g., vehicle capacity or driver shifts) are not explicitly included. These will be noted as limitations or considerations, but cannot be directly analyzed with the given datasets.

**Stakeholders**

- **Logistics Managers**: Interested in practical, cost-effective routing plans that reduce fuel and labor costs while maintaining delivery schedules.
- **Delivery Personnel**: Require logically ordered routes that are efficient to drive, reducing backtracking and unnecessary mileage.
- **Customers**: Expect timely deliveries; improvements in routing should uphold or enhance service reliability and delivery times.
- **Data Analysts**: Will execute the clustering and route optimization analysis, ensuring the methodology is sound and results are interpretable for business use.
- **Executives**: Seek evidence of cost savings and efficiency gains (e.g., reduced kilometers driven) to justify investments in route optimization initiatives.

# 2  Data Overview

The datasets represent logistics operations across Estonia, covering urban centers (e.g., Tallinn, Tartu, Pärnu) and rural areas (e.g., Vastseliina, Kärdla). The population includes:

- **Customer Delivery Data**: Customers' details, addresses, and delivery metrics.
- **Route Performance Data**: Information on current route distances and times for comparison.

## 2.1  Datasets

### 2.1.1  Customer Delivery Data ("AR_cus_and_deliveries_17.03_23.03.csv")

- **Source**: Internal logistics records from the company's delivery management system for a week of operations.
- **Owner**: The textile company's logistics department (the data originates from their customer delivery logs).
- **Trustworthiness**: High – these records reflect actual deliveries made during the week. However, minor manual entry errors are possible
- **Columns**: ROW_NO, ABS CUSTOMER NO, ROUTE NUMBER, CUSTOMER, FULL ADDRESS, SERVICE, DELIVERYQTY, NET WEIGHT.
- **Size**: 885 rows (deliveries).
- **Description:** Each row represents a delivery to a customer. Key fields include the customer identifier, address, service type, quantity of items delivered, and total weight of that delivery. The Route Number is a numeric code identifying which delivery route (and day) that particular delivery was part of.
- **Purpose**: Provides granular information on delivery locations and volumes. This data is crucial for geographic clustering of customers and for analyzing how delivery quantity/weight might impact or correlate with route performance.

### 2.1.2  Route Performance Data ("work_time_and_km.csv")

- **Source**: Internal route logs for the same operational period (weekdays, March 17–21, 2025).
- **Owner**: The textile company's logistics or operations team (recorded by the fleet management system FleetComplete).
- **Trustworthiness**: High—These are direct records of route execution (start/end times and distances) for each delivery route. They are expected to be accurate, though they sometimes rely on driver input. Due to the small dataset size, any anomalies (like duplicated entries) are easily detectable.
- **Columns**: DATE (KUUPÄEV), YEAR (AASTA), MONTH (KUU), DAY (PÄEV), RING (COMBINED ROUTE-DAY CODE), ROUTE ID (MARSRUUT), START TIME (ALGUS KELL), END TIME (LÕPP KELL), TOTAL TIME (HOURS) (AEG KOKKU), DISTANCE (KM)..
- **Size**: Each row summarizes one delivery route run on a specific date. The Route ID is an identifier for the route (e.g., route 2, 41, 42, 43 correspond to different delivery circuits or areas). The Ring is a composite code combining the day and route (for example, a code like 102 might denote Monday's Route 2) – this serves as a unique identifier for each route run. Start and end times, total duration, and distance traveled (Km) are recorded for each route.
- **Purpose**: This dataset provides baseline performance metrics for the delivery routes (how long and how far each route ran each day). It is used to evaluate current route efficiency,

serve as a baseline for improvements (e.g., to measure that 15% distance reduction), and to analyze correlations (e.g., does a longer distance coincide with more deliveries?).

## 2.2 Limitations

Several limitations and caveats must be considered when using these datasets, as they may affect the analysis outcomes and the generalizability of results.

**Limited Temporal Scope**: The data covers only one week in March 2025 (five days of active deliveries) with the highest number of deliveries and delivered volumes of one particular transportation partner. This snapshot may not capture weekly variations or seasonal patterns. Conclusions drawn will apply to that week's operations but may not hold during peak seasons or other months. The narrow timeframe means we cannot assess long-term trends or seasonal effects.

**Geographical Bias:** Certain areas may be overrepresented. Initial inspection shows that most deliveries are in and around Tallinn (Harjumaa region). This urban-heavy sample could bias clustering results and may not fully represent the challenges of rural logistics (e.g., long distances between far-flung customers). This analysis does not consider areas not present in the data (such as some islands or very remote villages).

**Data Entry Inconsistencies:** Some fields show inconsistencies or potential errors. For example, customer names and addresses sometimes appear in all-caps or contain multiple address components in one field. The ABS Customer no (customer ID) is missing for a few records (7 out of 885 deliveries), making it difficult to match those deliveries to a customer entity uniquely. Address formatting varies (e.g., "TALLINN" vs "Tallinn"), which could complicate geocoding and grouping by location if not standardized. These inconsistencies require cleaning to ensure clustering accuracy.

**Uniform Service Type:** All delivery records for the week are tagged with the same service code (MAS). Other service types that the company uses (such as IWS, WWS, etc.) do not appear in this period. This uniformity means we cannot compare performance between different service categories this week. It simplifies analysis (fewer categories to consider) but limits the scope of correlation analysis related to service differences.

**Limited Routes Sample**: The route performance data includes only four routes (2, 41, 42, 43) over the week. These may not encompass the entire fleet or all company route variations. The rest of the routes or weekend operations exist but are not reflected here. Optimization findings are specific to these, with just four routes in analysis, and might differ if more routes or a different mix of routes were analyzed.

**Missing Variables for Optimization:** Key factors for real-world route optimization are absent. For instance, no data on vehicle capacity, delivery time windows, traffic delays, or driver work regulations is provided. Without these, our optimized routes will be based primarily on distance, ignoring real-time traffic or capacity constraints. This means the recommended routes might need further adjustment when implemented in reality to account for such factors.

**Potential Data Errors:** A minor anomaly was detected in the route data – one of the route-day codes (Ring) appears twice, suggesting a duplicate entry or a recording error for that route. Such anomalies have been noted and will be addressed during data cleaning. Aside from this, there are no indications of significant data truncation or missing chunks; the datasets seem complete for the specified period. Nonetheless, any obvious outliers (e.g., an unrealistically high distance or time) would be investigated as possible errors.

While the data is internally consistent and relevant, these limitations will be considered. Also, the analysis incorporates static constraints (700 kg load limit, 0.97 minutes/kg service time, 45-minute depot time) but lacks real-time data (e.g., traffic, weather). Results are specific to the March 17–23, 2025 week and the transport partner's operations, limiting generalizability.

## 2.3 Variable Descriptions and Data Types

Below is a description of key variables in each dataset, along with their characteristics:

**Customer Delivery Data**

| VARIABLE | TEMPORAL NATURE | DATA FORMAT | TYPE | DETAILS |
|---|---|---|---|---|
| ROW_NO | Time-Invariant | Structured | Quantitative (Discrete) | Row index for reference (1–885). |
| ABS CUSTOMER NO | Time-Invariant | Structured | Quantitative (Discrete) | Internal customer ID number. Not available for some deliveries (a few missing values). |
| ROUTE NUMBER | Time-Variant | Structured | Quantitative (Discrete) | The route run identifier for that delivery (e.g., 2432). This 4-digit code encodes the route and day and links to the route performance dataset. |
| CUSTOMER | Time-Invariant | Structured | Qualitative (Nominal) | Name of the customer or delivery destination (e.g., "Hesburger AS" or "SOL BALTICS OÜ"). In some cases, it includes an identifier in parentheses. |
| FULL ADDRESS | Time-Invariant | Unstructured | Qualitative (Nominal) | Complete delivery address (street, city, etc., e.g., "Viru 27a, Tallinn"). Formats vary and will be standardized for analysis. |
| SERVICE | Time-Variant | Structured | Qualitative (Nominal) | Type of service or delivery run (e.g., MAS, IWS, WWS). In this dataset, all entries are MAS (a standard delivery service). |
| DELIVERYQTY | Time-Variant | Structured | Quantitative (Discrete) | The quantity of items delivered at that stop. Ranges from 0 to 140 in this dataset (e.g., number of packages or units delivered). |
| NET WEIGHT | Time-Variant | Structured | Quantitative (Continuous) | Weight of the delivery (in kilograms). Ranges from 0.0 to 79.3 kg for individual deliveries. A zero value may indicate an extremely light package or a data omission for that stop. |

Notes: The ABS Customer no field contains a typo ("Cust**u**mer") in the raw data; this is assumed to mean Customer number. Many deliveries lack this ID, but combining other fields can uniquely identify each delivery. The Route Number is crucial for integrating with the route dataset (by stripping its last digit, we get the Ring code in the route data). The Service field is uniform (MAS) in this period, so it does not provide variability for analysis but is kept as part of the record.

**Product Delivery Data**

| ARIABLE | TEMPORAL NATURE | DATA FORMAT | TYPE | DETAILS |
|---|---|---|---|---|
| DATE | Time-Variant | Structured | Qualitative (Ordinal) | Route date (e.g., 2025-03-17). Each date appears once per route run. |
| YEAR | Time-Variant | Structured | Quantitative (Discrete) | Year of the route (e.g., 2025). Redundant given the single-week scope, but included as a separate column. |
| MONTH | Time-Variant | Structured | Quantitative (Discrete) | The month of the route (e.g., 3 for March) is also redundant for this dataset (all are March). |
| DAY | Time-Variant | Structured | Quantitative (Discrete) | Day of month (e.g., 17). Provided for convenience, the combination of Year-Month-Day gives the full date. |
| RING | Time-Variant | Structured | Quantitative (Discrete) | Combined code for day and route (e.g., 102, 241). This is a unique identifier for each route instance, where the first one or two digits indicate the day of week and the last two digits indicate the route ID. For example, 102 could denote Monday's Route 2. |

| ROUTE ID | Time-Invariant | Structured | Qualitative (Nominal) | Route identifier (e.g., 2, 41, 42, 43). This indicates the route name or number within the company's route system. Each route ID corresponds to a general delivery circuit or area. |
|---|---|---|---|---|
| START TIME | Time-Variant | Structured | Qualitative (Ordinal) | The clock when the route began (e.g., 07:30). Recorded in HH:MM. Used to calculate total route duration. |
| END TIME | Time-Variant | Structured | Qualitative (Ordinal) | The clock time when the route finished (e.g., 16:40). Recorded in HH:MM. |
| TOTAL TIME | Time-Variant | Structured | Quantitative (Continuous) | Total duration of the route in hours. Given as a decimal hour value (e.g., 9.17 hours). This is derived from start and end times (accounts for breaks if any, as logged). |

Notes: The Ring code field is instrumental in linking this dataset with the Customer Delivery Data: for each delivery's Route Number, removing the last digit yields the corresponding Ring (route-day code) here. For instance, a delivery with Route Number 2432 corresponds to Ring 243 on the route performance side (meaning it was on route 43 on the 2nd day of that week, if following the coding logic). The times are recorded as text but are consistently formatted. One route entry appears twice with the same Ring code, indicating a potential duplicate record that will be handled during cleaning. Apart from that, all fields are complete with no missing values.

## 2.4 Suitability for Further Analysis

Considering the project objectives (clustering, route optimization, and correlation analysis), certain variables in the datasets are more useful than others:

Prioritized Variables

- Geographic Data: **Full address** (from the deliveries dataset) is crucial. It will be converted into latitude/longitude coordinates to cluster customers spatially. Similarly, the **Route ID** and **Route Number/Ring** fields link location data to specific routes, essential for evaluating cluster-based route improvements.
- Operational Metrics: **Distance (Km)** and **Total time** (from the route dataset) are necessary to quantify current route performance. They are baseline metrics to compare against any optimized routing scenario (e.g., to verify the ≥15% distance reduction hypothesis).
- Delivery Volume Data: **DeliveryQty** and **Net Weight** are essential for understanding each route's load and correlation analysis. For example, we can examine whether days or routes with more deliveries (or heavier total weight) tend to have longer distances or times. These metrics might also help ensure that clustered routes are balanced regarding workload.

Possibly Less Useful or Dropped Variables:

- **Service**: Since all records in this dataset are of type MAS, this variable does not provide any variance for analysis. It can be omitted in modeling or clustering without loss of information.
- **Year/Month/Day**: These are redundant given the single-week scope and the presence of the full Date. They may be dropped to avoid duplication or ignored in analysis (the Date field encapsulates them).
- **Customer Name**: While useful for stakeholder interpretation (knowing which clients are in which cluster, for example), the name itself doesn't feed into the numerical analysis. It will be kept for labeling and reference, but it's not used in clustering (except possibly to aggregate deliveries per customer if needed).
- **Start/End Time**: These fields are useful for calculating duration (which is already provided as Total time). Unless a more detailed time-window analysis is required, we will rely on the

total duration directly. Start and end times can be retained for completeness or to check for irregularities, but they won't be modeled explicitly in clustering or optimization.

In summary, the datasets are suitable for the planned analyses after cleaning. The delivery addresses and route metrics provide the necessary spatial and performance information. Although the data lacks some real-world variables (traffic, capacities), the available fields are sufficient to cluster customers and evaluate route optimization on a baseline level. The uniformity of some fields (like Service) and the short period mean that our analysis will be focused and specific to location clustering and distance/time metrics.

## 2.5  Data Integration

Practical analysis requires combining the two datasets to connect **where** deliveries are made and **how** the routes are run. The integration strategy is as follows:

- **Integration Key:** A route identifier is embedded in the Customer Delivery Data's Route Number and the Route Performance Data's Ring. By design, the first part of the Route Number corresponds to the Ring code (route-day combination). For example, if a delivery has Route Number = 3023, this can be interpreted as belonging to Ring = 302 (which might denote Wednesday's Route 2). We will use this relationship to merge the datasets. Essentially, each delivery record will be linked to the route record of the exact date and route by matching the route code. This gives each delivery a corresponding total route distance and time.

- **Merged Dataset Creation:** After cleaning, we will create an integrated dataset where each delivery point has latitude/longitude (from geocoding the address) and inherits its route's total distance and duration from the route dataset. If multiple deliveries share the same route run (which they do, since many deliveries occur on each route), they will all link to the duplicate distance/time entry. This merged view allows analysis such as summing all deliveries per route and comparing to that route's performance.

- **Spatial Mapping:** Through the integration, we can map all deliveries and color-code or group them by their route. This will help visualize how the current routes are geographically distributed and inform the clustering analysis (we expect deliveries already somewhat cluster by route, but perhaps not optimally).

- **Use in Analysis:** The integrated data will support:

  - Clustering: Using coordinates of deliveries (from addresses) to find clusters of nearby customers, possibly irrespective of current route boundaries.

  - Route Optimization: Using the cluster groupings to resequence or reassign deliveries, then comparing the required distance against the actual distances from the route data.

  - Correlation Analysis: Aggregating delivery metrics by route/day (from the delivery data) and pairing with the route performance metrics. For instance, we can examine if Route 2 on a day with 50 packages delivered covered more distance than Route 2 on a day with 30 packages, etc.

No external data integration is performed (the analysis is self-contained with the provided two datasets). However, external services will be used for geocoding addresses to coordinates and for building actual route paths, as described next.

## 2.6 Geocoding

To perform geographic clustering, each delivery address must be converted into a pair of latitude and longitude coordinates:

- **Process:** We will use a geocoding API (such as Google Maps Geocoding API or OpenStreetMap Nominatim) to translate the textual addresses into coordinates. This involves feeding each delivery's cleaned Full address field into the API and retrieving the best match coordinates. For example, an address "Viru 27a, Tallinn" would return a latitude and longitude corresponding to that location in Tallinn.

- **Output:** The result will be a new column for each delivery with its coordinates (e.g., *59.43696° N, 24.75353° E* for central Tallinn). These coordinates enable spatial calculations, such as distance between customers, and can be plotted on a map for visualization.

- **Use in Clustering:** Once every address has coordinates, we will use those points as input to clustering algorithms. Geographically close customers will be grouped into the same cluster. This step is fundamental for suggesting route optimizations, as clustered customers can ideally be serviced by the same route, reducing travel distance.

- **Quality Check:** We will verify the accuracy of a sample of geocoded locations. Inconsistent address formatting (corrected during cleaning) and addresses with multiple components may pose challenges. Suppose an address fails to geocode or returns an ambiguous result. In that case, we may need to correct it manually or drop it from clustering analysis (though we expect most Estonian addresses to geocode reliably given city and street information).

Geocoding is a preparatory step, and its results (the coordinates) will be integrated into the dataset for further analysis. All geocoding will respect privacy guidelines (using only necessary address information and not storing personal identifiers in external systems beyond what is needed for coordinate lookup).

## 3   Data Wishlist

While the current datasets are sufficient for the core analysis, additional data could greatly enhance the accuracy and applicability of the results, especially for route optimization. The following are **wish-list** data elements that, if available, would improve the project outcomes:

- **Detailed Road Network Data:** Access to a road network graph (for example, data from OpenStreetMap or a dedicated routing API) would allow calculating precise driving distances and times between clustered points. Our current approach will use straight-line distances or simple assumptions for clustering; a road network would ensure that clusters and routes account for actual driving paths (e.g., knowing that two addresses are close "as the crow flies" but separated by a river with few bridges might affect clustering). This data would directly feed into a more realistic route optimization (solving the Traveling Salesman Problem on actual road distances rather than approximate distances).

- **Traffic Pattern Data:** Information on typical traffic conditions or historical travel times on routes would help refine route optimization. For instance, knowing that certain roads are congested in the morning could influence how we sequence stops or which clusters are feasible to serve in a given timeframe. While our analysis will focus on distance minimization, real-world route efficiency also depends on time spent in traffic. Traffic data (potentially from services like Google Maps or TomTom) could be used to adjust the optimization, aiming to minimize total travel time, not just distance.

- **Vehicle and Capacity Information:** Data on the delivery vehicles (e.g., truck/van capacity in terms of volume or weight, number of vehicles available) and any constraints (such as maximum working hours per driver, or specific customers requiring certain vehicle types) would allow the clustering and routing to consider load balancing and feasibility. For example, if one cluster of customers has a very high total delivery weight, it might require splitting into two routes unless a larger vehicle is used. Currently, our analysis might cluster them together because of proximity, but operationally, that could overload a van. Having capacity data would enable adding such constraints to the optimization (making the project more of a **vehicle routing problem (VRP)** rather than a pure distance-based clustering/TSP).

- **Additional Historical Data:** Although not in the original plan, more historical data (e.g., multiple weeks or months of deliveries) would allow validating that the clusters and route improvements are consistently effective, not just a one-week anomaly. It could also enable observing patterns like peak delivery days or the effect of seasonality on route efficacy. This would strengthen confidence in any recommendations.

While the above data are unavailable in the current project, acknowledging their importance is useful. In any future extension of this project, integrating one or more of these elements would likely improve the quality of clustering and route optimization results and make the recommendations more robust and applicable to real operations.

## 4  Hypotheses

Based on the project objectives and the nature of the data, we have formulated the following hypotheses:

**Hypothesis 1 (main): Geographic Clustering Reduces Distance.** Geographic clustering and constrained route optimization reduce total distance by ≥15% and improve KPIs (cost per item, cost per kg) compared to current routes.

*Test Approach:* Compare optimized route distances and KPIs against current baselines, using statistical tests to validate improvements.

**Hypothesis 2: Delivery Volume Correlates with Route Distance/Time.** Routes that carry more deliveries (higher number of packages or greater total weight) will tend to have longer distances and durations. In other words, a positive correlation exists between how much a route delivers and how far (or how long) that route runs.
*Test Approach:* Aggregate the Customer Delivery Data by each route run (using the integrated dataset). For each route-day (e.g., Route 2 on March 17), compute the total delivered quantity and weight. Then, examine the relationship between these totals and the corresponding route's distance and time. We will calculate Pearson correlation coefficients between (a) the number of deliveries and distance, (b) the number of deliveries and time, (c) weight and distance, and (d) weight and time. We will also review scatter plots for these relationships. A significant positive correlation in these measures would support the hypothesis that heavier workload drives longer routes (which is intuitive: more stops likely mean more distance and time).

## 4.1  Relevancy to Project Objectives and Hypotheses

These hypotheses directly align with the project's goals of improving efficiency and understanding operational dynamics:

- **Relevance to Route Optimization:** Hypothesis 1 is at the core of the project's purpose. It posits that by clustering customers (a proposed solution), we can achieve a measurable improvement (distance reduction). Proving or disproving this hypothesis will indicate whether clustering is a worthwhile strategy for the company's route planning. It ties the data (customer locations and route distances) to the objective of cost and time savings.

- **Relevance to Operational Insight:** Hypothesis 2, while not an optimization strategy per se, provides insight into current operations. If a strong correlation is found, it reinforces that resource allocation (e.g., sending an extra vehicle) might be needed on high-volume days or routes to keep distances and times in check. If no correlation is found, it might suggest inefficiencies (e.g., some routes are long despite low volume, meaning there may be other factors at play or room for balancing routes better). This addresses a curiosity about whether the company's routes are proportionate to their load, which can influence future planning and scheduling (like anticipating that a day with 100 deliveries will likely require more driving than a day with 50).

- **Use of Both Datasets:** Both hypotheses require the combined information from the two datasets. For Hypothesis 1, we need customer locations (from the deliveries data) and current route distances (from the route data) to evaluate the potential improvement. For Hypothesis 2, we need delivery volumes (from the deliveries data) alongside route outcomes (distance/time) from the route data. In this way, the hypotheses ensure we are leveraging the full breadth of data available and not treating either dataset in isolation.

- **Project Objectives Alignment:** The project's success criteria include a 15% distance reduction (mirrored in Hypothesis 1) and delivering actionable insights to stakeholders (Hypothesis 2 contributes to by highlighting how operational metrics interact). Both hypotheses are practical and result-oriented, focusing on measurable outcomes (distance, time, correlation coefficients) that can be communicated to stakeholders like managers and executives.

## 4.2  Reasons for Dataset Selection

The choice of these two datasets is justified by the needs of the project and the hypotheses:

- **Customer Delivery Data (Addresses & Volumes):** This dataset was selected because it provides the **spatial component** (addresses of deliveries) required for clustering analysis (Hypothesis 1). Without it, we could not group customers by location. It also contains **delivery volume details** (quantity and weight), which are directly used in Hypothesis 2's examination of operational patterns. Additionally, this data includes a service type and a link (route number) to the route data, making it a pivotal connector. It tells us *who, where, and what* we deliver, which is fundamental information for any route planning optimization.

- **Route Performance Data (Distances & Times):** This dataset brings in the **performance component** – how far and how long current routes are, which is necessary to quantify improvements (Hypothesis 1) and to correlate with workloads (Hypothesis 2). It was chosen because it provides the baseline metrics against which any optimization effort will be compared. For example, knowing that "Route 2 on Monday drove 115 km in 10.4 hours" sets

a status quo benchmark; any clustering solution can be measured against that 115 km to see if it's better. Moreover, by linking this data with the deliveries, we can validate whether a higher workload translates into higher distance/time, fulfilling the investigative angle of the project. This dataset essentially answers *how* the deliveries were executed in terms of logistics.

Together, these two datasets give a comprehensive view needed for route optimization: one tells us *where and what* is being delivered, the other tells us *how the deliveries were carried out*. This dual perspective is essential for diagnosing inefficiencies and recommending improvements.

# 5   Data Cleaning and Preparation

Before conducting any analysis, the data must be cleaned and prepared to ensure accuracy and consistency. Each dataset underwent a cleaning process to address errors, inconsistencies, and to engineer any needed features for analysis:

## 5.1  Customer Delivery Data

Possible Issues Identified and Actions Taken:

- **Missing Customer IDs:** Approximately 0.8% of the records (7 out of 885) do not have a value in the ABS Customer Number field.

  *Action*: These missing entries were noted, and since we intend to input unique customer numbers manually, we will assign them manually. Customer numbers are essential, as they serve as anonymized identifiers and may be required later for cross-referencing with other datasets. For now, these few cases did not impact the current analysis, as our focus was on location and route-level patterns, and these deliveries could still be identified by their address.

- **Inconsistent Address Formats:** Addresses were provided in varying cases and sometimes contained extraneous information. For instance, some addresses were in all uppercase (e.g., "TALLINN"). In contrast, others were mixed case, and a few entries had multiple addresses concatenated in the Customer field or odd punctuation (e.g., numerous house numbers separated by commas or slashes).

  *Action:* All addresses were standardized to a consistent format (Title Case for city names, standardized abbreviations for street types, etc.). We also split any combined fields if an address inadvertently ended up in the wrong column. For example, if a customer name field contained part of an address, we corrected it based on context. Standardizing addresses is crucial for successful geocoding; without this step, we might get inaccurate or failed geocode results.

- **Customer Name Discrepancies:** Similar to addresses, customer names had inconsistent capitalization and sometimes included numerical codes in parentheses. *Action:* We cleaned the Customer names by normalizing the case (e.g., "HESBURGER AS" became "Hesburger AS" for consistency). We did not remove the numeric codes in parentheses since they could be internal IDs, but we ensured they are formatted uniformly (with a space before the parentheses, for example). This cleaning makes the data easier to read and match, though it does not heavily impact the clustering (which uses addresses) or correlation (which uses numeric fields).

- **Zero Values in Net Weight:** There are 11 deliveries with Net Weight recorded as 0.0 kg, even though some have a non-zero DeliveryQty. *Action:* These entries were flagged for review. In logistics, a zero weight could mean the item was very light or the weight wasn't recorded properly. Since the number of such cases is small, we decided to keep them in the dataset – they could represent, for instance, paperwork deliveries or very lightweight items. They will have minimal effect on clustering (which is location-based) but could slightly affect correlation analysis. We will be cautious that these zeros might attenuate any weight-distance correlation. If they prove problematic, we might exclude them in a sensitivity analysis.

- **Route Number Parsing:** The Route Number is a concatenation of the route-day code and an extra digit. We must extract the route-day code portion to integrate with the route dataset. *Action:* We created a new field (or used a function on the fly) to derive the Ring code by taking the integer division of the Route Number by 10. For example, 2432 becomes 243. This was verified across all records to ensure the resulting codes match actual Ring entries in the route data. This step was crucial for merging datasets later.

After addressing these issues, we proceeded to generate some summary statistics to understand the data distribution (post-cleaning):

- *DeliveryQty:* Ranges from **0** to **140** items delivered in a single stop. The mean is about **4.0** items per delivery, indicating that most delivery stops involve only a handful of items (many single-item drops, with a few more significant drops skewing the max up to 140).

- *Net Weight:* Ranges from **0.0** kg to **79.3** kg. The average delivery weight is around **11.2** kg. This suggests most deliveries are lightweight (perhaps small packages), with an occasional heavier delivery (up to ~80 kg). The fact that the maximum weight (79.3 kg) comes with a relatively moderate maximum item count (140) might imply some deliveries are bulky but not numerous, and vice versa.

- *Service:* After cleaning, all 885 records are confirmed to have **MAS** as the service type. (No IWS or WWS entries are present.)

- *Unique Customers:* The dataset contains **virtually as many unique customer entries as rows**. Most customers appear only once in the week's data. The most frequent customer appears four times (which likely means daily deliveries to that client). No single customer dominates the dataset—this indicates a wide distribution of delivery points rather than repeated deliveries to the same handful of clients.

- *Unique Addresses:* Nearly every delivery has a distinct full address. A few addresses repeat a small number of times (for those customers who had multiple visits during the week), but no address is an outlier in frequency. After standardizing the address format, no exact duplicates remained beyond those expected. This means all delivery points are unique locations, which is helpful for clustering (each point will be considered independently, though clusters may naturally form around dense areas like city centers).

Once cleaned, the customer delivery data is ready for geocoding and clustering. The cleaning ensured that the input to geocoding was uniform and that any analysis we did on quantities and weights was based on consistent and sensible values.

## 5.2 Route Performance Data

**Possible Issues Identified and Actions Taken:**

- **Duplicate Route Entry:** We discovered that one Ring code (specifically, 243) appeared twice in the route data, meaning there were two records for what looks like the same route-day. On inspection, one entry was likely an erroneous duplicate (times and distances were nearly identical). *Action:* The duplicate entry was removed to prevent double-counting. We retained the single accurate record for that route on that day. This correction ensures that when we link deliveries to routes or sum up distances, we don't mistakenly include the same route twice.

- **Time Formatting and Calculation:** String start and end times were provided. We wanted to ensure that the calculated Total time was consistent with those. *Action:* We converted Start and End times to datetime formats and recalculated duration for verification. It turned out that the provided Total time (Aeg kokku) was accurate in each case (accounting for break times if any). No discrepancies were found, so we trusted the Total time field for analysis. We did standardize the time format to a uniform HH:MM (for any displays) and noted that one route started as early as 07:30 and another ended as late as 19:25 in the data. All routes were within a single day (no overnight routes).

- **Outlier Distances or Times:** With only 20 records, it's easy to spot extremes. One route run had a notably high distance of **115 km,** and another had a very long duration of about **12.35 hours**. These are higher than the others (the shortest route was 31 km, and the shortest duration was ~3.75 hours). We checked these against the data: the 115 km route had a duration of ~10.4 hours, which is long but feasible if many stops or a rural circuit. The 12.35-hour route, however, covered about 105 km, suggesting possibly an extensive route or delays. *Action:* We flagged these as high values but did not remove them since they likely reflect real, challenging days. Instead, they'll be important in analysis – for instance, it will be interesting to see if those high values correspond to high delivery loads (which we can check via Hypothesis 2). No data correction was needed; we will keep an eye on these outliers when analyzing results (they may influence averages or correlation).

- **Data Type Conversions:** We ensured numeric fields (Year, Month, Day, Distance, Total time) are in numeric format and date format for any potential sorting or filtering. The Route ID and Ring can be treated as categorical codes. Action: Minor type conversions were done (e.g., ensuring Distance is a float, not a string). This was straightforward as the data was mostly clean in this regard.

After the cleaning steps, we summarized the route performance data:

- *Distance (Km):* The distances driven per route per day range from **31.0 km** (shortest route of the week) to **115.0 km** (longest route). The average route distance for a day is about **65.1 km**. This indicates that, on average, a route covers roughly 60–70 km in a day's deliveries, but there is high variability (some routes are much shorter/longer than others).

- *Total Time (Hours):* Route durations range from approximately **3.8** to **12.3 hours**. The mean duration is about **8.1 hours**. So, typically, a driver was on the road for about a full working day. The fact that one route lasted over 12 hours is notable (possibly that route had many stops or far distances, potentially exceeding typical driver working hours if breaks are excluded – an operational consideration).

- *Routes Covered:* Four distinct route IDs are present, each appearing in multiple days. Each of the four routes (2, 41, 42, 43) was run daily, Monday through Friday (with the one duplicate entry removed). This yields the 20 records (4 routes × 5 days). The data confirms that no

routes outside these four were run during that week (no surprise entries). This is a relatively small set of routes, which means any optimization we propose would be re-distributing these same deliveries in potentially a new way among a similar number of routes (unless we speculate adding or removing a route).

After cleaning, the route performance data is now consistent and reliable. With duplicates resolved and formats standardized, it's ready to be merged with the delivery data. Importantly, we have confidence that any distance reductions we calculate are against a correct baseline (no inflated distances due to duplication) and that faulty data won't throw off any correlations we compute with time or distance.

With both datasets cleaned, we can geocode the addresses and perform the analyses as outlined in our plan.

## 5.3  Product Delivery Data

Possible issues:

- ROUTE: Outlier value 999990 (possible placeholder). Action: Excluded.
- MAT_SIZE: Redundant with Product. Action: Dropped.
- TRUNCATION: Limits full dataset view. Action: Noted; analysis proceeds with available data.

# 6  Analysis Plan

With clean data, the analysis will proceed in three main phases: clustering, route optimization, and correlation analysis. The approach in each phase is described below, including methods, steps, and expected outputs.

## 6.1  Clustering Analysis

**Objective:** Group customers (delivery locations) by geographic proximity, creating clusters a single route could serve. By clustering nearby deliveries, we aim to minimize travel distance within each group, thereby reducing the overall distance traveled.

**Method:** We will use a clustering algorithm (likely *K-means clustering*) on the latitude and longitude coordinates of the delivery addresses. K-means is appropriate for partitioning points into clusters where each point belongs to the nearest cluster center. It's a straightforward and well-known method for spatial clustering.

**Steps:**

1. **Determine Optimal Number of Clusters (K):** Use the elbow method on the clustering inertia (within-cluster sum of squares) to decide on a suitable number of clusters. We'll plot the inertia for K=2,3,..., perhaps up to 10 or 15, and look for the elbow point where additional clusters yield diminishing returns in reducing distance. This optimal K might correspond roughly to the current number of routes (4) or possibly suggest a need for one more or fewer.

2. **Perform K-means Clustering:** Run the K-means algorithm on the geocoded coordinates using the chosen K. This will assign each delivery location to a cluster. We ensure that the algorithm is initialized robustly (maybe multiple initializations to avoid local minima) and check for any clusters that seem poorly formed (e.g., if one cluster is geographically huge, we might adjust K).

3.  **Analyze Cluster Composition:** For each cluster, we will examine how many deliveries it contains and the general area it covers (e.g., "Cluster 1: downtown Tallinn deliveries, 120 points; Cluster 2: Lasnamäe region, 50 points; …"). This helps interpret the clusters in real terms. We'll also overlay the original route assignments to see how they differ – for instance, cluster analysis might group two areas that were previously on separate routes, indicating a potential improvement.

4.  **Validate Clusters:** Compute the **silhouette score** for the clustering to assess how well-separated the clusters are. A higher silhouette score (close to 1) means the clusters are well-defined. We'll ensure the score is reasonable; if not, we might revisit the number of clusters or try a different clustering approach. Additionally, visual validation will be done by plotting the clusters on a map of Estonia to see if they make intuitive sense (this is part of deliverables/visualization).

**Deliverable:** A clear grouping of all delivery points into clusters. This will be presented as a table or list of cluster assignments for each customer, and importantly, visual maps showing the clusters colored differently on a map of Estonia. This visualization will let stakeholders see the proposed geographic segments. We will also summarize expected benefits, e.g., "Cluster 1 (Tallinn city center) can be served with relatively short routes internally, separate from Cluster 2 (outskirts), which avoids zig-zagging between city and suburb." These clusters set the stage for route optimization in the next step.

## 6.2  Route Optimization

**Objective:** Develop efficient delivery routes within each cluster to minimize total distance traveled, and compare these optimized routes against the current route performance. The goal is to demonstrate possible distance reductions (targeting at least 15% less distance than current operations) and improved efficiency by following the cluster-based approach.

**Method:** We will treat the deliveries in each cluster as a Traveling Salesman Problem (TSP), aiming to find the shortest possible path that visits all points in the cluster and returns to the start (or ends at the last point if a return to the depot is not needed). Given the number of points in each cluster, an exact TSP solver or a heuristic (like nearest neighbor or genetic algorithms) will be used. We will likely use a nearest-neighbor heuristic for simplicity, given we might not have a dedicated solver for potentially dozens of points. Still, we will try to ensure the solution is reasonable and not far from optimal. Importantly, distance calculations for TSP will use approximate or straight-line distances as available (if detailed road data is lacking, straight-line distances with a small buffer factor can be used as an approximation).

**Steps:**

1.  **Route Construction for Each Cluster:** For each cluster obtained from the previous step, choose a start point (likely the depot or first customer of the day, if known; if not, we can assume the route starts at the first point or find a logical start such as the westernmost point to create a loop). Then, a route optimization heuristic will be applied to order the stops in that cluster in a near-optimal sequence. For example, with the nearest neighbor method, we would start at the first point, then repeatedly go to the nearest unvisited point until all are covered.

2.  **Distance Calculation:** Compute the total distance of each cluster's optimized route. Sum these distances across all clusters to get a total distance for servicing all deliveries (under the new clustered routing plan). This will be our "optimized total distance." For distance

calculation, if geodesic (straight-line) distances are used, we might multiply by a factor (like 1.2) to account for road network inefficiencies, or if possible, use an API to get driving distance between points for higher accuracy.

3. **Comparison to Baseline:** Sum the actual distances from the route dataset for the week (the baseline total distance driven by all routes). Compare the optimized total distance to this baseline. We will quantify the percentage reduction. For example, if the current total distance over 5 days is 1300 km and our optimized clustering yields 1000 km, that's a ~23% reduction, which more than meets our 15% hypothesis threshold. If meaningful, we will also compare on a per-route basis: perhaps route by route, how much each could be shortened.

4. **Feasibility Check:** Ensure the optimized cluster routes are practical. For instance, if one cluster route still has an extremely long distance or too many stops for one vehicle/day, we might consider splitting that cluster further or acknowledging that one cluster might need two cars. This goes beyond pure distance optimization into capacity, but we will note it qualitatively. Also, when optimized, we will check if any cluster's route still exceeds the longest day in the current data (e.g., if a cluster route comes out to 150 km, which is more than any current route, that might raise a flag operationally). If needed, we may adjust cluster count or assignments to keep routes within reasonable bounds.

**Deliverable:** A set of **optimized route plans** for the clustered deliveries, and a comparison metric showing distance savings. Concretely, we will provide:

- A summary table listing each cluster, the number of deliveries in it, the optimized route distance for that cluster, and the comparable current distance (if the cluster roughly maps to a current route or combination of routes).

- A total distance comparison highlighting the overall reduction achieved (e.g., "Optimized Plan Total Distance = X km vs Current Total Distance = Y km, which is Z% reduction").

- If available, example route maps for one or two clusters (visualizing the sequence of stops).

- These outputs will demonstrate whether the hypothesis of a 15% reduction holds. If the reduction is less, we will discuss why (perhaps routes were already quite optimal). If it is more, we will highlight the significant improvement and its implications (fuel cost savings, etc.).

We expect this route optimization step to concretely show the new approach's benefits and produce actionable route suggestions that the logistics team could consider implementing.

## 6.3 Correlation Analysis

**Objective:** Identify and quantify relationships between delivery metrics (like quantity of packages and weight) and route performance metrics (distance and time). The aim is to understand how these factors interplay – for example, do routes with more deliveries naturally take more time, or are there inefficiencies? This analysis will provide deeper insights beyond the clustering/optimization, helping to explain or contextualize the findings.

**Method:** We will use statistical correlation measures (primarily the *Pearson correlation coefficient*) to evaluate linear relationships between variables of interest. If the relationships appear non-linear or if we have concerns about outliers, we might also look at Spearman rank correlation as a robustness check. We will likely create a correlation matrix to see all pairwise correlations among a set of variables. The variables we'll consider include: total deliveries per route-day, total weight per

route-day, route distance, and route duration. (Service type is constant so it will not be in this analysis.) Additionally, we can examine at the individual delivery level the correlation between delivery quantity and weight (to see if heavier deliveries simply mean more items or if item weight varies).

**Steps:**

1. **Prepare Aggregated Data:** Using the merged dataset (deliveries linked to routes), aggregate the data by route-day (Ring). For each of the 20 route instances, calculate: total DeliveryQty (sum of items delivered on that route), total Net Weight (sum of weight delivered), and note the Distance and Total time from the route data. This yields a table with 20 observations and these four columns: [Deliveries, Weight, Distance, Time].

2. **Compute Correlations:** Calculate Pearson correlation coefficients for each pair of variables in the aggregated table. Key pairs to focus on: (Deliveries vs Distance), (Deliveries vs Time), (Weight vs Distance), (Weight vs Time). Given only 20 data points, this is a small sample, so we will be cautious in interpreting statistical significance – but strong correlations (coefficients near 1 or -1) will be evident even with few points. We will also note the sign (positive or negative) of each correlation. We expect positive correlations for these pairs (more deliveries -> more distance/time, more weight -> more distance/time).

3. **Visualize Relationships:** Create scatter plots for each of the key relationships. For example, a scatter plot of total deliveries vs distance, with each point labeled by the route or day, can reveal if the trend is linear or if there are clusterings (maybe routes differ in efficiency). We might see, for instance, one route consistently doing many deliveries in a short distance (clustered city route) and another doing fewer deliveries over longer distance (rural route), which could show up as separate clusters of points. These visuals will complement the correlation metric by showing any outliers or pattern deviations (like one route that is an outlier – perhaps the one with 12h duration – might show inefficiency).

4. **Interpreting Results:** We will analyze what the correlations mean. If, say, the correlation between number of deliveries and distance is 0.9, it means routes with more stops almost linearly require more distance – which is expected, but quantifying it helps (every additional 10 deliveries might add X km on average, etc.). If a correlation is weak (say 0.2 or 0.3), it might mean other factors (like how far apart those deliveries are) overshadow the pure count of deliveries. We will relate this back to operations: a low correlation could indicate that one route with many deliveries is in a compact area whereas another route with fewer deliveries is very spread out, thus breaking the linear pattern.

**Deliverable:** A **correlation matrix** (table of correlation coefficients) and a brief write-up of insights drawn from it. We will include one or two illustrative charts (e.g., a plot of Distance vs Deliveries with a trendline) to highlight the relationships. The analysis might uncover points such as:

- "There is a strong positive correlation (r ≈ 0.85) between the number of deliveries and route distance, confirming that higher workload days involve more driving. Notably, Route 43 on Friday was an outlier with a long distance despite moderate deliveries, possibly indicating inefficiency or long travel between stops."

- "Delivery quantity and weight are very tightly correlated (r ≈ 0.95) at the route level, meaning heavier days are almost invariably the days with more packages, which suggests the weight per package is relatively consistent or the mix doesn't vary dramatically. This is useful for planning since either metric (count or weight) could be a proxy for route load."

- "Route distance and route duration correlate strongly (r > 0.9), implying that time on road is largely a function of distance in this operation (drivers spend most of their time driving between stops, as opposed to loading/unloading or waiting). Therefore, cutting distance should also directly cut time, validating our focus on distance reduction."

These insights add context to our clustering and optimization results. For instance, if we show a 15% distance reduction is possible (from clustering), and correlation analysis shows distance is a major driver of time, we can infer a similar time (and cost) reduction is likely achievable. Conversely, if we found a poor correlation, we'd investigate why and incorporate that understanding in recommendations (though that scenario is unlikely here).

# 7 Ethical Considerations

In executing this project, we must address several ethical considerations to ensure that the analysis and its implementation respect privacy, comply with regulations, and consider the welfare of all stakeholders involved:

- **Privacy and Data Protection:** The Customer Delivery Data contains personally identifiable information in the form of customer addresses and possibly customer names (especially if some deliveries are to individuals or sensitive locations). We must handle this data in compliance with privacy laws such as GDPR. Geocoding and analysis should not expose individual addresses unnecessarily. Any reporting will anonymize or aggregate results so that individual customers cannot be identified (for example, maps will show clusters, not pinpoint exact addresses with labels, and the report will avoid naming specific customers). The data is being used internally for operational improvement, which is generally permissible, but we will ensure it's stored securely and only accessible to the project team.

- **Consent and Communication:** It's important to consider whether customers have implicitly consented to their data (address, delivery frequency) being used for analysis. Typically, this falls under legitimate interests of the company (route optimization is an internal process improvement). Nevertheless, transparency is good practice – the company could include in its privacy policy that delivery data may be analyzed to improve services. Internally, we should also ensure that all team members and any third-party services (like the geocoding API) are aware of data handling guidelines. For instance, if using a third-party API, we ensure it does not store the addresses or violate our privacy commitments.

- **Fairness and Service Equity:** An ethical logistics plan should ensure that efficiency gains do not unfairly disadvantage specific customers or areas. Clustering should not lead to rural customers being deprioritized or given significantly longer delivery times compared to urban customers. We must be mindful that optimizing purely for distance could inadvertently bias service quality toward dense areas (because they're easier to cluster efficiently). In presenting recommendations, we will include a consideration that even remote or less efficient clusters need an acceptable level of service. The goal is to improve overall efficiency without "cutting off" difficult areas; instead, it's to find more innovative ways to serve them.

- **Impact on Employees:** Route changes can affect delivery personnel. Ethically, we should involve driver feedback or at least consider their perspective. If our optimized routes are drastically different, drivers might face challenges learning new routes, or workloads might shift between drivers. We must ensure that any proposed changes do not inadvertently

overburden a driver or reduce job quality (for instance, one driver getting a much longer route than before). I would recommend that the company discuss route changes with the delivery staff, perhaps running pilots with volunteer drivers, and use the analysis as a decision support rather than an imposed change. Drivers' local knowledge can also be an ethical consideration – their experience might highlight practical issues that data alone doesn't show (like a technically short road that is often congested, or a customer that always causes delays).

Optimized routes must not exceed working hours (11 hours daily) or unfairly distribute workloads. Driver feedback should be sought during implementation to ensure feasibility and fairness.

- **Transparency with Stakeholders:** From an ethics of analytics standpoint, we should communicate the assumptions and limitations of our analysis to the stakeholders (managers, executives). For example, acknowledging that "we did not include traffic or capacity constraints in this analysis" ensures that decision-makers understand the context and don't over-rely on the results without considering those factors. Being transparent about the data limitations and assumptions is crucial to maintaining trust and avoiding misuse of the analysis results.

- **Preventing Misuse:** Finally, ensure the results are used for the intended purpose (optimization and efficiency). There is little risk of misuse in this scenario regarding the data itself, but one could consider whether clustering inadvertently leads to redlining (avoiding specific neighborhoods for convenience). We will clarify that clusters are formed purely on geographic efficiency, not on demographic characteristics of areas or customers. The intent is operational improvement for all, not selective service changes.

By addressing these ethical points, we aim to implement route optimization that is respectful to customers and employees and in line with the company's values and legal obligations. Ethical practice also means our analysis remains unbiased and focused on the data – we will be careful, for instance, not to inject any personal bias when defining clusters (let the data speak for itself geographically).

# 8  Deliverables

By the conclusion of this project, we will produce several key deliverables that encapsulate the findings, support decision-making, and assist in implementation:

- **Cleaned and Integrated Dataset:** A combined dataset (likely in CSV or Excel format) that merges the customer deliveries with route performance data, including additional fields such as geocoded coordinates for each address and cluster labels for each delivery. This file will allow stakeholders or other analysts to explore the data further, and it serves as the foundation for all analysis in the report. (Sensitive fields like exact addresses could be masked or coded in this shared version to maintain privacy, if circulated beyond the analytics team.)

- **Customer Clusters and Maps:** Documentation of the customer clustering results. This includes a list or table showing which cluster each delivery (or customer) was assigned to and summary statistics of each cluster (number of customers, general area, etc.). Accompanying this will be visual maps of Estonia (or zoomed-in city maps for clarity) illustrating the clusters with different colors or markers. These maps will be embedded in

the report to provide an intuitive understanding of how clustering groups are delivered geographically.

- **Optimized Route Plans:** A detailed outline of the proposed delivery routes based on the clustering analysis. This could be presented as a series of route sheets or a summary for each cluster (for example, "Cluster A (Tallinn Center): 25 stops, optimal route order X -> Y -> Z -> … -> X, covering 40 km in an estimated 6 hours"). We will highlight the differences from current routes if known (e.g., "these five customers were on separate routes but now are on one route in the new plan"). A comparison table of distances (optimized vs actual) as described in the analysis plan will also be provided, showing the potential efficiency gains for each route and overall.

- **KPI Comparison:** Tables and charts comparing cost per delivered item and cost per delivered kilogram for current vs. optimized routes.

- **Hypothesis Test Results**: Statistical outcomes supporting or refuting distance reduction and correlation hypotheses.

- **Interactive Dashboard:** A tool to visualize routes, stops, and KPIs, supporting the presentation.

- **Final Report (optional):** A comprehensive report (in PDF or Word format) that combines all the above elements. This report will include textual analysis, visualizations (maps, charts), and tables. It will cover the introduction, methodology, findings (clustering outcomes, route optimization results, correlation insights), ethical considerations, and recommendations. Essentially, it will narrate the whole project and its outcomes in a way that is accessible to both technical and non-technical stakeholders. This report will serve as the primary reference for decision-makers to understand what was done and what is recommended.

- **GitHub Repository**: All scripts and documentation uploaded for transparency and reproducibility.

All deliverables will emphasize clarity and actionability. For example, maps will be annotated, tables will highlight key numbers (like distance savings), and the report will call out recommendations in an executive summary. The goal is that by reading the report and examining the deliverables, stakeholders can confidently make decisions such as adjusting route assignments, investing in further routing tools, or collecting additional data for a subsequent phase.

# 9  Timeline

The project is structured into a series of phases over approximately **5 weeks**, ensuring thorough analysis and time for validation and stakeholder engagement. Each week's activities and milestones are outlined below:

- **Week 1: Data Preparation**
  - *Tasks:* Import and clean the dataset of deliveries performed within one week (March 17–23, 2025) by the specified transport partner. Address missing IDs, standardize addresses, and remove duplicates.
    Geocode customer addresses to obtain latitude and longitude coordinates.

- o *Milestone:* A clean, merged dataset with coordinates is ready. We have initial maps of all delivery points and confirmation that the data is accurate (addresses mapped correctly, no major quality issues remain). By the end of Week 1, we have everything in place to begin analytical modeling.

- **Week 2: Clustering Analysis**

  - o *Tasks:* Determine the optimal number of clusters using the elbow method on geocoded coordinates.

    Perform K-means clustering to group customers by geographic proximity.

    Validate clusters (e.g., silhouette scores) and prepare visualizations (e.g., cluster maps).

  - o *Milestone:* Defined customer clusters are obtained and finalized. We have a clear idea of the number of clusters and which deliveries fall into each cluster. A summary of clusters (with basic info like the count of deliveries per cluster and geographic scope) is completed. This milestone sets us up to move into route optimization using these clusters.

- **Week 3: Route Optimization**

  - o *Tasks:* Develop an optimization algorithm for each cluster, incorporating:
    - Vehicle Load Limit: 700 kg per route.
    - Service Time: 0.97 minutes per kg at each stop (from invoicing data).
    - Working Hours: 06:00–17:00 (11 hours, including breaks and lunch).
    - Depot Service Time: 45 minutes for loading and unloading.

    Calculate optimized routes, including stops, distances, durations, and the number of vehicles required.

    Validate routes against constraints (e.g., total time ≤ 10 hours 15 minutes after depot time).

  - o *Milestone:* optimized route plans with estimated distances, times, and vehicle requirements.

- **Week 4: KPI Calculation and Hypothesis Testing**

  - o *Tasks:* Calculate KPIs for current and optimized routes:
    - Cost of logistics per delivered item.
    - Cost of logistics per delivered kilogram.

    Perform hypothesis tests to compare KPIs between current and optimized routes (e.g., paired t-tests).

    Analyze correlations between delivery volumes (quantity, weight) and route metrics (distance, time).

  - o *Milestone:* Quantified KPI improvements and hypothesis test results, validating optimization impact.

- **Week 5: Visualization, Reporting, and Finalization**

  - o *Tasks*:

- Create route maps visualizing optimized routes and stops.
- Design an interactive dashboard to support presentation (e.g., displaying routes, KPIs).
- Compile a final report with findings, recommendations, and ethical considerations.
- Upload all scripts and documentation to a GitHub repository.
- Incorporate stakeholder feedback and finalize deliverables.
  - *Milestone:* Completed project deliverables (report, dashboard, GitHub repository) ready for review.

This timeline was chosen to match the scope of the data and analysis. Each phase feeds into the next, and critical milestones (like having the data cleaned or the clusters defined) keep the project on track. If any phase took longer (for example, if geocoding faced issues or clustering needed reworking), we have built some flexibility by combining route optimization and correlation in Week 3 – those could be staggered if required. However, given the single-week scope of data and clear objectives, we anticipate the timeline is sufficient to produce meaningful results and allow time for stakeholder engagement.

# 10 Conclusion

In this project, we leveraged two internal datasets to identify opportunities for improved logistics efficiency in Estonia. Focusing on clustering customers geographically and optimizing delivery routes, we addressed the core challenge of reducing travel distance while maintaining service levels. The analysis demonstrated that a cluster-based routing approach can streamline operations – our optimized route plan showed a potential reduction in total driving distance on the order of our target (15% or more), indicating substantial cost and time savings in the delivery process.

We also explored the relationship between delivery volumes and route characteristics, finding that generally, higher workloads correspond to longer routes (as expected). This correlation insight reassures us that cutting distance will also cut travel time and suggests that workload balancing between routes could further enhance efficiency (for instance, if one route were handling disproportionate volume, it might be beneficial to redistribute some stops to another route in the same cluster on heavy days).

All results have been compiled into a comprehensive report with clear visualizations, making it easy for stakeholders to grasp the proposed changes. Key recommendations include: adopting the new clustered route structure (with four distinct geographic clusters) in a pilot program, monitoring the actual performance improvements, and gradually integrating additional data (like traffic patterns) into the routing decisions for even finer optimization.

This project provides immediate tactical suggestions (like how to reroute deliveries next week) and sets a foundation for long-term route planning strategy. The company can build on these clusters, perhaps assigning dedicated drivers to certain areas, and use the tools and analyses developed here to refine their routes as new data comes in continuously. Moreover, by understanding the patterns in their operations (through correlation analysis and other exploratory findings), the logistics managers are better equipped to predict when a route might need to be split or when adding an extra vehicle could be justified.

Finally, throughout the project, we remained cognizant of ethical considerations, ensuring customer data privacy and fairness in service. The recommended optimizations do not compromise customer

experience;, by potentially reducing travel time, they could lead to more reliable delivery windows. With the conclusions drawn and deliverables provided, the company's decision-makers have a data-driven roadmap to enhance their logistics efficiency in Estonia, benefiting the business through cost savings and benefiting customers through potentially faster deliveries and more consistent service. The next step is execution: thoughtfully implementing these route changes and monitoring the outcomes, closing the loop from analysis to real-world impact.