

1. To validate the given hypotheses, the analysis of the data must proceed to prove that such an approach is worth considering. With the current knowledge, the validation of the hypothesis should be based on a few steps:

A. Firstly, the current data should be split into two sets: training and test. The split should be randomized, but it should also maintain the distribution of mean delivery times within each sector.

B. Afterwards, based on the naïve and sector-based models, predictions are made on the training data to obtain results.

C. Then, validation of the results should be performed using simple metrics like mean absolute error (MAE) or root mean square error (RMSE).

D. The results might also be visualized on charts for easier interpretation. Based on the results obtained from the validation and visualization process, it should be easy to determine whether the hypothesis holds true.

E. If it's true, then it should be implemented further and used for the whole dataset. Otherwise, the company should stick to the naïve method.

2. An easy method to improve accuracy is by implementing a regression model, which can capture the linear relationship between input features and the target variable (delivery time). It's easy to implement and a good starting point for further improvements. The methodology to validate the new method would be the same as described above: splitting the data into training and test sets and comparing the results with the naïve method.

Alternatively, with a larger amount of data, it might be worth considering implementing neural networks, particularly recurrent neural networks, which can capture dependencies in sequential data. This would be useful for predicting delivery times. The methodology for validating this method would be the same as described in the first point.

3. Deliveries could take more time due to various reasons. For example, as mentioned in the task, some buildings do not have elevators and might have narrow staircases. Additionally, some buildings might have security measures that significantly extend delivery time, such as requiring visitors to sign in or obtain special badges. Customer availability could also be an issue; some deliveries require the customer to be present, so if the recipient is not available, the courier may need to make multiple delivery attempts, which extends delivery time for others. High demand or peak periods, such as holidays or promotional events, can result in an increased volume of deliveries, which may extend the duration of delivery.

4. Data that might be worth considering gathering would include, for example, seasonal trends and analysis of delivery demand and order volumes. Such data can help determine the expected delivery time during specific time periods. It is also worth considering gathering data related to weather conditions in which packages are delivered, as difficult conditions might significantly increase delivery times. Additionally, tracking cost and revenue data to monitor costs associated with delivery

operations, such as fuel, labor, or vehicle maintenance, is not only essential for delivery time optimization but also for cost management of such services.

5. The risk of overestimating or underestimating delivery times is quite similar in both cases, as they both can result in wasting money and resources. When overestimating, a company wastes resources such as drivers and vehicles, leading to increased operational costs and decreased efficiency. It may also lead to customer dissatisfaction since they receive deliveries later than expected. In the second case, when underestimating, it may result in rushed operations, leading to increased risk of errors, accidents, and safety issues. Like the previous case, it may also lead to customer dissatisfaction since the customer doesn't receive a package within the expected timeframe.