




# An analysis of the impact of demand aggregation on the solution quality for facility location problems

## *Análise do impacto do agrupamento da demanda na qualidade da solução de problemas de localização de instalações*

Renata Akemi Marçal Imai<sup>1</sup>, Claudio Barbieri da Cunha<sup>1</sup>, Cauê Sauter Guazzelli<sup>2</sup>

<sup>1</sup>Universidade de São Paulo, São Paulo, São Paulo, Brasil

<sup>2</sup>Inteligência de Negócios e Pesquisa Operacional, São Paulo, São Paulo, Brasil

Contact: renata-imai@alumni.usp.br,  (RAMI); cbcunha@usp.br,  (CBC); caue@inpo.eng.br,  (CSG)

### Submitted:

16 August, 2022

### Accepted for publication:

16 October, 2023

### Published:

20 December, 2023

### Associate Editor:

Renato da Silva Lima, Universidade Federal de Itajubá, Brazil

### Keywords:

Demand aggregation.  
Facility location.  
K-means.

### Palavras-chave:

Agregação da demanda.  
Localização de instalações.  
K-means.

DOI: 10.58922/transportes.v31i3.2801

### ABSTRACT

Inspired by real-world applications, this paper studies the impact on the quality of solutions for facility location problems in which demand points are aggregated to reduce the size of the underlying mathematical formulation. Two aggregation methods are analyzed and compared: demand points aggregated based on municipal boundaries or other similar administrative boundaries as usually done in practice and using the K-means clustering algorithm. Regarding a business-to-business (B2B) distribution context, two datasets comprising the location of thousands of drugstores in Brazil were generated, and 18 different instances of the fixed cost facility location problem were derived. The results show that solutions with aggregated demand points by municipality yield a maximum 0.43% difference in the objective function value in comparison with the respective disaggregated mode, while the difference using K-means algorithm did not exceed 0.03%. We also performed an in-depth analysis of the regions where the demand points were allocated to distinct selected facilities in the aggregated and disaggregated models. It was possible to observe that in the model with aggregated demand points by municipality, differences in transportation costs are greater than using the K-means clustering algorithm as the aggregation procedure. This suggests that aggregating demand points with the K-means clustering algorithm yields both better objective function values, and selected facilities closer to demand points in the cases where the resulting assignment of demand points to the selected facilities is not the same as the results of the unaggregated model.

### RESUMO

Inspirado em problemas reais, neste trabalho é avaliado o impacto na qualidade da solução de problemas de localização de instalações nas quais os pontos de demanda são agregados a fim de reduzir o tamanho do modelo matemático resultante. Dois métodos de agregação são analisados e comparados: pontos agregados considerando os limites geográficos dos municípios ou quaisquer outras subdivisões administrativas como feito na prática, ou obtidos com a aplicação do algoritmo *K-means*. Considerando um contexto de distribuição *business-to-business* (B2B), foram gerados dois conjuntos de dados contendo a localização de milhares de farmácias no Brasil, a partir do qual 18 instâncias do problema de localização de instalações não capacitadas foram derivadas. Os resultados indicam que soluções com agrupamento de pontos de demanda por município levam a uma diferença de, no máximo, 0,43% no valor da função objetivo em comparação ao modelo desagregado correspondente, enquanto essa diferença não excedeu 0,03% quando utilizado o algoritmo *K-means*. Também foi realizada uma análise das regiões nas quais os pontos de demanda dos modelos agregado e desagregado foram alocados a diferentes instalações. Foi possível observar que nos modelos que consideram pontos de demanda agregados por município, as diferenças nos custos de transportes são maiores que na aplicação do algoritmo *K-means*. Isso indica que agrupar os pontos de demanda utilizando o algoritmo *K-means* proporciona valores melhores de função objetivo, bem como instalações mais próximas aos pontos de demanda nos casos em que os resultados são diferentes aos do modelo desagregado.



## 1. INTRODUCTION

A facility location problem comprises selecting the number and locations for facilities, which can be non-capacitated or capacitated, to serve a set of demand points or customers

in the best possible way. Facility location problems are among the most fundamental problems in combinatorial optimization (Laporte, Nickel and Saldanha-da-Gama, 2019). It is a well-established research field that has been very active since the 1960s (Saldanha-da-Gama, 2022) due mainly to its application in many real situations. It has become increasingly relevant in the context of designing supply chains and distribution networks that comprise locating facilities such as manufacturing plants, distribution centers, and warehouses such that the demands of all customers are met while respecting the capacities of the facilities and the total cost, comprising fixed and variable transportation components, is minimized.

Distinct facility location problems, algorithms, and solution methods have been studied extensively in the literature and can differ in how they are classified concerning location space: continuous or discrete. In the former, facilities can be positioned anywhere in some continuous regions, while in the latter one has to select the best location for facilities from a given set of potential candidates. According to Saldanha-da-Gama (2022), discrete facility location problems are predominant in logistics and transportation. Exact and heuristic methods can be used to determine the best solution for discrete location problems. Both capacitated and uncapacitated facility location problems are combinatorial optimization problems classified as NP-hard, which makes the computation of optimal solutions in reasonable times unlikely (Wolsey, 1998; Fischetti, Ljubić and Sinnl, 2016; De Armas et al., 2017). Thus, heuristics and metaheuristics are more commonly utilized and even necessary when the size of the problem, given by the number of candidate facilities and demand points is elevated, as well as the additional complexity brought about by specific constraints that do not allow exact methods (i.e., integer programming-based mathematical formulations) to be employed.

In location modeling, two types of aggregation are commonly used: aggregation of continuous data into discrete points and aggregation of a large dataset of discrete points into a smaller one (Daskin et al., 1989). As pointed out by Francis et al. (2009), location problems occurring in urban or regional settings may involve a large number of demand points, whose aggregation is a common way to obtain tractable models. In some cases, it may be impossible, and unnecessary, to include every demand point in the corresponding model. Also, spatial aggregation allows decreasing data collection and modeling costs, but at the expense of not working with the actual locations, which may affect mainly the transportation costs.

The perceived need to handle individual points in strategic location decisions within the context of logistics and urban distribution can sometimes lead to unnecessary complexities and challenges. While the intention to account for every unique demand point is valid, it is important to weigh this against the practical implications and computational limitations. The strategic level of locational decisions encompasses inherent inaccuracies of various kinds, primarily revolving around forecasting future demands and costs. Incorporating spatial aggregation techniques can be a pragmatic way to strike a balance between accuracy and feasibility in strategic location decisions. By thoughtfully grouping demand points, decision-makers can attain efficient, actionable insights without getting entangled in the unnecessary intricacies posed by dealing with individual points and all the underlying calculations.

In this paper, we investigate how methods for spatially aggregating demand points can influence the quality of exact solutions when modeling and solving to optimality the fixed

charge facility location problem (UFLP). This exploration becomes especially pertinent due to the inherent challenges encountered in tackling larger real-world instances that arise in practice. The expected direct consequence is the reduction of model size in a way to allow optimal solutions to be obtained that otherwise would either take much longer running and oftentimes excessive times or even could not be solved at all. In other words, by implementing aggregation strategies, we aim to mitigate these obstacles, resulting in a streamlined model, enabling the derivation of optimal solutions that would otherwise remain out of reach due to computational limitations posed by factors such as model size or the required solving time.

Our approach stems from the fact that in many real-world distribution systems, it is common to observe regions in space where demand points, such as commercial outlets, residences, etc., exhibit some sort of spatial concentration as urban forms evolve. One typical and commonly used way to aggregate demand points is to consider official or administrative boundaries such as city, district, neighborhood, and county limits, as it usually requires none or very little effort to be obtained as such attributes are present in databases used for such purposes. Also, geographical locations (i.e., latitude and longitude for every address) are also becoming widely available.

However, one question that arises is whether they yield a good representation of the problem, once a municipality, for instance, may have areas with a different number of demand points or distinct amounts of product units to be delivered, so that representing all the demand in a point in the center of the area could misrepresent the real problem, especially in terms of the distances that are directly related to the resulting total variable cost of distribution. As highlighted by Janjevic, Winkenbach and Merchán (2019), clustering-based methods typically begin by grouping customers together and subsequently utilize these clusters as a substitute for location decisions. Yet, we have not come across literature that comprehensively examines this empirically prevalent spatial aggregation approach on a city or municipality level from a scientific standpoint.

In this context, this paper proposes to analyze two clustering methods and compare the results of each aggregated model. We then analyze the cases in which the aggregated model is different from the unaggregated one and identify regions in space where misallocated demand points are. Our motivation is to identify if the K-means heuristic (MacQueen, 1967), a well-known and established clustering algorithm, can be employed as a reliable clustering method, capable of reducing total processing time to solve the mathematical model without affecting the overall quality of the solutions, as we change the baseline scenario in the fixed-charge facility location problem. To the best of our knowledge, this is the first time in the literature that such a comparison between the results of a classical geographical clustering with an algorithm-based clustering method, such as the K-means is presented.

The remainder of this paper is organized as follows: Section 2 presents a brief literature review on exact methods for a class of facility location problems, clustering demand points, and errors associated with demand point spatial aggregation. Section 3 describes the proposed approach, while Section 4 details the two sets of demand points and the respective instances generated. Section 5 presents the computational experiments. Lastly, Section 6 presents the conclusions of this work.

## 2. LITERATURE REVIEW

Facility location problems, algorithms, and solution methods have been studied extensively in the literature. Saldanha-da-Gama (2022) provides a recent review of the role of facility location in logistics and transportation, highlighting not only the challenges brought about by a fast economy globalization and a strong increase in environmental concerns but also current trends and future challenges.

In our research, we selected three main topics that are worth addressing in this review: (i) facility location problems, (ii) demand point spatial aggregation methods, and (iii) demand point aggregation errors. For this literature review, each one of them is briefly covered in this section.

Melo, Nickel and Saldanha-da-Gama (2009) highlight that facility location is a well-established research area and has been extensively studied in the scientific literature. The facility location problem consists of finding the best location for a set of candidates, satisfying the constraints to attend to the demand, minimizing the total distance, time, or cost. As some classes of facility location problems have already been meticulously explored in scientific papers (ReVelle, Eiselt and Daskin, 2008), this literature review will focus only on the fixed charge problems, hereafter represented by the abbreviation FCFL. Daskin (2013) provides a comprehensive review of conventional problems from which other more complex and realistic models are derived.

Fixed charge problems are used as test models in optimization models since they assume that facilities have a capacity large enough to supply all the demand, and the number of open facilities is endogenously defined (Hakli and Ortacay, 2019). Initially, to solve the FCFL, only exact methods were employed, and they aimed to obtain optimal solutions to the problems. Some works on exact methods of FCFL used branch-and-bound (Khumawala, 1972) and other variations such as branch-and-cut (Aardal, 1998), the lagrangian relaxation (Barcelo et al., 1990), the dual problem (Erlenkotter, 1978), or simply linear programming (Van Roy, 1986), for instance. Studies on the polyhedral structure of the problem's convex hull were conducted by Aardal, Pochet and Wolsey (1995). A recent paper related to modeling and the analysis of polyhedral structures was written by Sankaran (2007).

Recent papers on the FCFL are employing heuristic methods to obtain optimal or near-optimal solutions, given the increase in the problem's size, since exact methods may not always find a solution to these problems. Regarding heuristic methods, a distinction must be made. There are problems solved only using heuristic methods, as Kratica et al. (2001) and Hakli and Ortacay (2019), but there are also problems that employ heuristics to improve some exact methods, such as Galvão and Raggi (1989), Barahona and Chudak (2005), and Hansen et al. (2007).

As highlighted by Goodchild (1979), many of the fields to which facility location models are employed take the database as some aggregation of a geographically dispersed demand. As the size of the problem increases, the total processing time also increases. One possibility to handle large facility location problems is through demand point aggregation, which consists of clustering two or more demand points into only one that is representative of the set. In literature, two aggregation methods are usually employed. The first one clusters the demand points in relation to their geographical location. For example, Zhao and Batta (1999) and Hodgson, Shmulevitz and Körkel (1997) aggregate the demand points based on the centroid of the area's zip code or census tract,

respectively. In turn, Sankaran (2007) suggests the aggregation of demand points in large metropolitan areas in two clusters: one that represents the largest municipality, and the other that represents the remaining municipalities.

While geographical clustering techniques are relatively simple to implement, they may not always accurately identify regions with varying demand densities in space. Therefore, an alternative approach is to utilize established clustering algorithms presented in the literature. For instance, O'Kelly (1992) employs the K-means algorithm to cluster demand points in a hub location problem, whereas Erkut and Bozkaya (1999) utilize both K-means and linkage methods to aggregate the demand points and compare their results in a p-median problem. Paul and MacDonald (2016) also employ K-means clustering to aggregate demand points and obtain solutions for a facility location problem by means of heuristics as solution methods.

In a distinct approach, Tong and Church (2012) analyzed the effects of aggregating continuous spatial units into discrete points within the context of the location set covering problems. The authors propose a measure to understand and quantify errors associated with a continuous aggregation scheme, and demonstrate its concepts with an empirical study of sitting emergency warning sirens.

In addition to the geographical clustering and the usage of clustering algorithms, some researchers, such as Levin and Ben-Israel (2004) and Gao (2021), employ heuristic approaches to determine the optimal configuration of  $K$  clusters within their datasets, while Andersson et al. (1998) propose aggregating demand points in grids. Recent papers, by Cebecauer and Buzna (2017), Irawan et al. (2017), and Irawan and Salhi (2015) aim to integrate clustering as part of the optimization process, using heuristics to obtain optimal or near-optimal solutions, and minimizing the errors from using aggregated data.

Errors can arise when using aggregated demand data to solve our facility location problems. One of the seminal papers related to this subject was written by Hillsman and Rhoda (1978). According to the authors, there are three sources of errors, named source A, B, and C. Later, Hodgson, Shmulevitz and Körkel (1997) proposed Source D errors, and Erkut and Bozkaya (1999) present errors that occur when handling problems with aggregated demand points. In several works related to calculating errors that arise from aggregation, authors calculate the gap between the optimal solution from the unaggregated demand point model and the optimal solution from the aggregated model (Andersson et al., 1998; Cebecauer and Buzna, 2017; Erkut and Bozkaya, 1999; Hodgson, Shmulevitz and Körkel, 1997; Irawan and Salhi, 2015; Zhao and Batta, 1999).

Francis et al. (2009) surveyed the literature for aggregation approaches to a large class of location models, comprising median, center, and cover problems and compared various aggregation error measures. Their focus, however, was different: to determine how many aggregated demand points are enough to provide an accurate representation in the case of home deliveries where each private residence might be a demand point. Their assumption is that in such cases it may be impossible, and unnecessary, to include every demand point in the corresponding model.

In its turn, Jacobs-Crisioni, Rietveld and Koomen (2014) examined the errors arising in spatial aggregation considering the scale and shapes of aggregated areal units, considering the Modifiable Areal Unit Problem (MAUP). The authors argue that when it comes to urban analysis, several spatial variables should be taken into consideration to avoid data loss

when aggregating the data; suggesting that using greater spatial units to aggregate demand points can lead to aggregation errors due to data loss.

In the literature review, we noticed that there are many methods to solve the FCFL and to aggregate demand points. Nevertheless, the solutions of aggregated models are mostly analyzed as a gap in the unaggregated model. It was also spotted that only Erkut and Bozkaya (1999) carried out a comparison between clustering algorithms, not geographical clustering method.

### 3. PROPOSED APPROACH

In this section, we present the formulation of the FCFL we use in the paper. We also show the proposed approach to aggregate the demand points; to calculate the gap between the analyzed models; and to identify locations where a demand point is allocated to different facilities in aggregated and unaggregated models.

Our decision to utilize the FCFL stems from the fact that utilizing uncapacitated facility location models in strategic decision analysis offers several distinct advantages that align with the nature and goals of this level of analysis. The absence of capacity constraints shifts the focus toward strategic considerations. This allows decision-makers to explore different scenarios and evaluate the impacts of various facility placements on a larger scale. These models facilitate a more comprehensive examination of alternative strategies and their potential outcomes. Furthermore, in the context of strategic decision analysis, the size of the facility assumes a relatively diminished significance. At this level of planning, where the primary focus is on long-term positioning and broader resource allocation, the intricacies of facility size become less critical. Instead, the emphasis rests on understanding optimal placement, overall network design, and macro-level impacts. This approach allows decision-makers to consider a range of facility sizes in the strategic context, without becoming overly entangled in the details that pertain to operational capacity. This broader perspective promotes more flexible and agile decision-making, enabling the exploration of different facility sizes in alignment with the strategic objectives of the organization.

The FCFL consists in allocating an undetermined number of facilities, minimizing the sum of fixed setup costs and variable costs of serving the demand points from these facilities (Verter, 2011). We chose the FCFL to handle strategic decisions since the facilities' sizes and capacities can be further adjusted without imposing restrictions on the model.

The problem that we formulate is deterministic, static, discrete, single-echelon, single-objective, where each demand point (vertex) can be serviced from one or more open facilities, and where no inventory decisions are relevant.

Let  $I$  be the set of candidate facilities to be located and  $J$  be the set of demand points (or nodes), which can be either the original ones or the ones obtained by some aggregation method whose location is known. We also define:

- $f_i$  as the fixed cost of opening and operating candidate facility  $i \in I$ ;
- $Q_j$  as the required amount of product (demand) of point  $j \in J$ , denoting either the unaggregated or aggregated amount, in which case it corresponds to the sum of the demands of all nodes that are represented by  $j \in J$ ;
- $c_{ij}$  as the unit transportation cost from candidate facility  $i \in I$  to demand point  $j \in J$ .

It is important to note that in the cases where  $j \in J$  denotes a subset of demand points that have been aggregated, the value of  $c_{ij}$  is computed regarding the geographical location of the corresponding node that represents those points.

The decision variables are:

- $Y_i$  is a binary variable that indicates whether a candidate facility  $i \in I$  is selected to open or not;
- $X_{ij}$  that indicates the total amount of units sent from a candidate facility  $i \in I$  to a demand point  $j \in J$ .

We assume that  $f_i > 0$  for all  $i \in I$ ;  $Q_j > 0$   $j \in J$ ; and  $c_{ij} \geq 0$  for all  $i \in I, j \in J$ .

Hence the problem can be formulated mathematically as follows (Daskin, 2013):

$$\text{minimize } \sum_{i \in I} \sum_{j \in J} c_{ij} \cdot X_{ij} + \sum_{i \in I} f_i \cdot Y_i \quad (1)$$

subject to:

$$\sum_{i \in I} X_{ij} \geq q_j, \forall j \in J \quad (2)$$

$$X_{ij} \leq q_j \cdot Y_i, \forall i \in I, \forall j \in J \quad (3)$$

$$Y_i \in \{0,1\}, \forall i \in I \quad (4)$$

$$X_{ij} \geq 0, \forall i \in I, \forall j \in J \quad (5)$$

The objective function (1) aims to minimize the sum of fixed and variable costs. The restrictions (2) ensure that the demand in each point  $j \in J$  is satisfied, while restrictions (3) ensure that only open facilities can supply the points demand. Restrictions (4) and (5) define the domain of the decision variables. To adapt the FCFL to solve problems with aggregated demand points, the subscripts  $j$  in the formulation presented in (1)-(5) are changed with subscripts  $k$ . With respect to the demand point clusters, each cluster  $k \in K$  consists of a subset of demand points  $j \in J$ , and the unit transportation cost represents the cost of attending the demand of clusters  $k \in K$  from facilities  $i \in I$ .

To reduce the size of the resulting formulation, we then apply two different aggregation methods to the demand points. Aggregation Method 1 (M1) clusters the demand points considering their geographical location in space; more specifically, we identify the respective municipalities where they are located (Sankaran, 2007). It means that all demand points within a municipality are represented by a single point (centroid) located at their gravity center. It should be highlighted that any other geographical administrative division or boundary such as county, borough, or district could be used instead, without loss of generality. Our decision to consider municipalities stems from the fact that, as in many real-world problems found in practice, our experiments consider a nationwide distribution system for a country of continental dimensions. In such scenarios, demand point aggregation is a commonly used practice to reduce the computational complexity of the problem as well as its underlying mathematical formulation.

The second method, called Aggregation Method 2 (M2), employs the K-means clustering algorithm (MacQueen, 1967) to aggregate demand points, without necessarily considering the boundaries of the municipalities to determine the clusters of the demand points. One of the advantages of using K-means algorithm is that it allows pre-setting the number of clusters. Therefore, in this paper, we fix the number of clusters,  $K$ , to be the equal to the number of municipalities in the datasets. This choice enables a fair comparison between different aggregation methods at the same level of granularity, thus yielding mathematical models of equivalent size in terms of number of decision variables and constraints. Such a comparison facilitates the evaluation of results and enhances their interpretability.

After aggregating the demand points, the respective facility location models are solved to optimality optimized. First, the unaggregated demand point model, in which the number of demand points is equal to the number of clusters  $K$ , is solved. This corresponds to the model described in Equations 1-5 with subscripts  $j$ . This model is solved so we can compare solutions with the aggregated models. Once the set of open candidates has been determined, the clusters that represent the aggregated demand points are dismantled (i.e., the aggregation is undone) and each point is allocated to the open facility with the lowest associated variable cost. This step ensures that any errors associated with source C are eliminated from the analysis.

To enhance the communication of the results of our analysis and facilitate the comparison of results between the unaggregated and aggregated demand points models, we introduce the following definitions:

- $m_A$  as the set of facilities selected in the aggregated facility location problem;
- $m_U$  as the set of facilities selected in the unaggregated facility location problem;
- $f(m_A, c_{ij})$  as the objective function calculated using  $m_A$  and the respective variable costs  $c_{ij}$  that takes into account aggregated demand points;
- $f(m_U, c_{ij})$  as the objective function calculated using  $m_U$  and the respective costs  $c_{ij}$  that takes into account unaggregated demand points.

The expression utilized to calculate the gap between aggregated and unaggregated solutions is noted in Equation 6. The gap value indicates whether there are differences between the aggregated and unaggregated models and if the aggregated model is over or underestimating the objective function compared with the unaggregated model.

$$gap = \left[ \frac{f(m_A, c_{ij}) - f(m_U, c_{ij})}{f(m_U, c_{ij})} \right] \quad (6)$$

Apart from comparing the gap between the optimal solutions of the unaggregated and aggregated demand point models, we also analyze the differences in allocation and transportation costs at the municipality level. A local analysis can also help identify regions in space where such misallocation occurs and associate it with higher or lower demand point density regions. One method to identify these regions is by applying the Density-Based Spatial Clustering of Applications with Noise (DBSCAN) algorithm (Ester et al., 1996), which clusters points based on the notion of density, i.e., how close to each other the demand points are in space. DBSCAN clustering enables us to identify possible metropolitan or megacities regions without any prior knowledge.



*Eps* is possibly the most important parameter to be set in DBSCAN and should be properly calculated (Han, Kamber and Pei, 2012). If it is too big, it can lead to fewer clusters than necessary or if it is too small, it possibly identifies too many clusters. We use Rahmah and Sitanggang (2016) method to calculate the *Eps* parameter. First, a distance matrix containing distances between every pair of points is determined. These distances are then ordered ascending, with respect to the third nearest neighbor, and the inflection point of this curve is then set as *Eps*.

When employing DBSCAN to cluster demand points, the following possible outcomes can arise: (i) all demand points in a municipality are clustered together; (ii) there are multiple clusters within the municipality; (iii) the cluster spans across multiple municipalities; or (iv) there is at least one cluster within a municipality boundary and at least one other cluster with demand points in more than one municipality (a combination of situations ii and iii).

#### 4. INSTANCE SETS

The two clustering approaches were evaluated using two demand point sets (designated Set 1 and Set 2) for the purpose of this research. Both demand sets are derived from a large nationwide real business-to-business (B2B) distribution system in Brazil.

It is worth noting that we opted to create our own datasets instead of relying on established well-known instances from the literature, such as OR-Lib (Beasley, 1990) and M\* (Kratika et al., 2001). The reasons for this choice are that (i) these instances lack the necessary geospatial attributes required to aggregate demand points by geographical region, and hence, prevent us from analyzing the effect of aggregation on regions with multiple demand points; (ii) these instances exhibit geometric features that are overly regular or uniform, which does not accurately reflect the irregular shapes and patterns typically observed in real-world scenarios.

Data scraping techniques were used to obtain the addresses and other attributes of the demand points, which correspond to stores requiring service from regional distribution centers whose number and locations should be determined with the aim of minimizing total fixed and variable transportation costs.

Table 1 presents the principal attributes of Set 1 and Set 2. Set 2 is a variant of Set 1, formed by random draws to select demand points initially chosen for Set 1. This approach was taken to alleviate potential bias present in the demand point initially chosen for Set 1. Notably, while the cumulative demand of Set 2 mirrors 80% of that in Set 1, the number of municipalities encompassed by Set 2 matches 90% of the municipalities found in Set 1. Across both Sets, the number and the placement of candidate locations are the same.

**Table 1:** Basic dimensions of input data for Set 1 and Set 2

	Set 1	Set 2
Demand points	5,574	4,460
Municipalities	650	585
Total demand (units)	633,265	498,872
Candidate locations	150	150

Table 2 presents the size of unaggregated and aggregated models for Sets 1 and 2. For both sets, the size of aggregated models M1 (clustering by municipality) and M2 (applying K-means clustering) is the same. To compare two different aggregation methods at the same aggregation level, we set the number of clusters to be equal to the number of municipalities existing in the Set. Hence, Models M1 and M2 allow a reduction of 8.57 times the number of lines and columns in Set 1, while in Set 2, the number of decision variables is reduced by 7.51 times.

**Table 2:** Resulting model sizes for Set 1 and Set 2

	Set 1		Set 2	
	Unaggregated Models M1 & M2		Unaggregated Models M1 & M2	
Size	150 x 5,547	150 x 650	150 x 4,460	150 x 585
Lines	841,674	98,150	673,460	88,335
Columns	836,250	97,650	669,150	87,900
Non-zero values	2,508,300	292,500	2,007,000	263,250

For each demand point set, nine different instances were generated. Each instance differs with respect to fixed costs. Based on the information provided by Bueno (2019) and Rodero (2018), we consider a baseline fixed cost of \$9,000. To analyze how problems with aggregated demand points behave when changing fixed costs, we increase the base value by 2, 5, and 10 times, but we also decrease it by 2, 5, 10, 100, and 1,000 times. The unitary transportation costs are based on a national transportation fare in effect in Brazil (Brasil, 2021). The instance label represents to which data set they belong, and which fixed cost they have. For example, instance 2-90 represents a problem with Set 2 demand points when the fixed cost is \$90.

## 5. COMPUTATIONAL EXPERIMENTS

In this section, we report the results of the computational experiments. We solve the FCFL for each instance set and their solutions are analyzed. All mathematical formulations were built using Python 3.8 and optimized using the software Gurobi version 9.1.1 (Gurobi, 2022) in an Intel®Core™i5-10210U @1.60 GHz computer with processor, 8 GB, 64-bit memory operating under Windows 10.

### 5.1. Evaluation of aggregation effect in instances of Set 1

Table 3 presents the solution for unaggregated and aggregated models of Set 1 instances. Regarding the objective function value, the aggregated models overestimate the unaggregated models in instances 1-9 to 1-1800. Model M1 is the one that yields the largest gap in comparison with the unaggregated model (0.43% in instance 1-9). As the fixed cost of opening a new facility increases, the objective function also increases.

With respect to the number of open facilities in the optimal solutions, instances 1-9 and 1-90 are the only ones in which aggregated models have fewer open facilities than the unaggregated model. However, it is worth noting that when considering a different and smaller set of open facilities, the reduction in fixed costs in these scenarios is less

apparent, primarily due to a subsequent rise in transportation costs. As a result, this disparity becomes more pronounced in the context of aggregated models.

In instances 1-900 and 1-1800, the difference in the objective function of aggregated models indicates that there is at least one open facility different from those opened in the solution of unaggregated model. Hence, for presenting a different set of open facilities, aggregated models present these gaps in objective function values. For instances 1-4500 to 1-90000, aggregated and unaggregated models yield the same objective function values as the same sets of open facilities are selected and the allocation of the aggregated and unaggregated demand points to them is also identical.

With a maximum acceptable gap of 0.50%, the solutions provided by the models utilizing aggregated demand points can be considered part of an alternative solution set for the unaggregated problem. Consequently, these alternatives should be weighted thoughtfully by the decision maker, who must strategically determine the best option over time.

**Table 3: Set 1 solutions**

Instance/Model	1-9			1-90			1-900		
	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2
Objective Function	69,870.58	70,171.78	69,891.68	75,778.89	75,958.03	75,780.89	107,286.45	107,286.69	107,316.10
Difference	-	0.43%	0.03%	-	0.24%	<b>0.00%</b>	-	<b>0.00%</b>	0.03%
# Open Facilities	122	107	119	83	80	82	37	37	37
Processing time	26	1	1	189	1	1	425	1	1
Explored nodes	168	0	0	1,294	0	0	384	0	0
Simplex iterations	8,855	173	206	27,925	547	716	48,533	2,219	2,452
Instance/Model	1-1800			1-4500			1-9000		
	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2
Objective Function	127,962.48	127,977.07	127,962.48	171,926.72	171,926.72	171,926.72	222,405.04	222,405.04	222,405.04
Difference	-	0.01%	0.00%	-	0.00%	0.00%	-	0.00%	0.00%
# Open Facilities	28	28	28	18	18	18	12	12	12
Processing time	162	1	1	181	1	1	378	1	1
Explored nodes	1,457	0	0	1,738	0	0	4,561	0	0
Simplex iterations	55,717	3,208	3,413	81,833	4,983	5,293	283,492	6,806	6,851
Instance/Model	1-18000			1-45000			1-90000		
	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2
Objective Function	299,313.99	299,313.99	299,313.99	456,279.80	456,279.80	456,279.80	580,209.97	580,209.97	580,209.97
Difference	-	0.00%	0.00%	-	0.00%	0.00%	-	0.00%	0.00%
# Open Facilities	9	9	9	5	5	5	4	4	4
Processing time	159	1	1	430	2	2	582	2	2
Explored nodes	1,604	0	0	2,998	0	0	3,910	0	0
Simplex iterations	156,286	8,616	8,948	373,587	13,192	13,643	506,408	17,357	17,623

The number of Gurobi's simplex iterations is also associated with problem complexity. On average, the unaggregated model needs 171,404 simplex iterations to obtain the optimal solution, while Models M1 and M2 require 6,344 and 6,571 iterations, respectively. As expected, the unaggregated model requires more iterations and therefore processing time to find the optimal solution, differently from the aggregated models. Processing time is 281 seconds (approximately 5 minutes) on average for the unaggregated model, and only one second on average for Models M1 and M2, which indicates that reducing the size of the problem can considerably reduce the processing times.

To better understand exactly how the two aggregation methods affect the optimal solutions in terms of selected facilities and respective demand points assigned to them, we employ DBSCAN algorithm (Ester et al., 1996) to cluster the data to identify the regions where differences between Models M1 and M2 and the unaggregated model arise. The *Eps* parameter

considered for this dataset was 0.020 and was calculated using the nearest neighbors' distance between demand points. The results show that DBSCAN was able to identify 15 municipalities with more than one cluster, 16 municipalities with demand points in clusters with other municipalities, and 13 municipalities with both clusters with their own demand points and with other municipalities. The remaining municipalities did not meet the criteria and were considered neither large municipalities nor belonging to metropolitan areas.

We then selected instance 1-9 of Set 1 to be analyzed, since it is the one that yields the largest difference between aggregated and unaggregated models. For each municipality, transportation costs with respect to unaggregated and aggregated models are calculated to be compared.

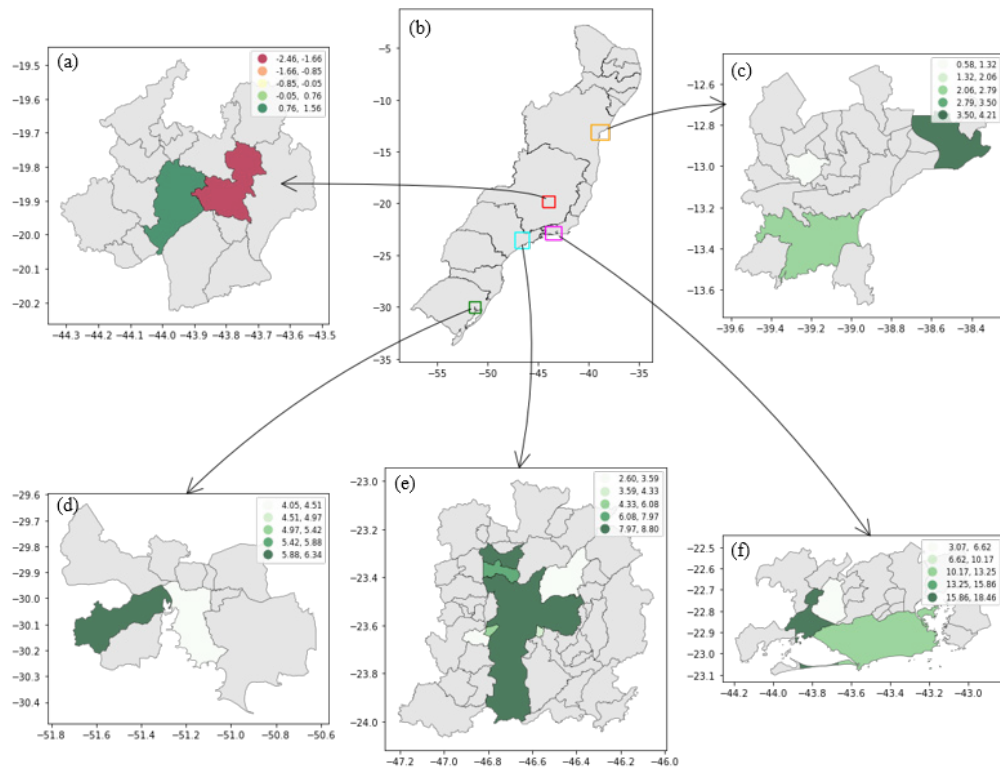
Figure 1 depicts the regions where Model M1 has some demand points allocated to different facilities in the aggregated and unaggregated models. The five regions shown in the figure are associated with municipalities previously identified by DBSCAN. Differences in transportation costs range between -2.46% and 18.46%, which means that there are municipalities where the transportation costs in the aggregated model underestimate or overestimate the real accurate cost resulting from the unaggregated model. This is expected since there are differences in open facilities in each model, especially due to demand point aggregation.

For Model M2, Figure 2 shows the locations of demand points allocated to different facilities in aggregated and unaggregated models and the differences in transportation costs. All highlighted areas were identified using DBSCAN algorithm and are common to Model M1. This suggests that some regions may be more prone to present misallocated points than others due to the number of demand points and the aggregation method employed. Differences in transportation costs range between -1.37% and 12.30%. This range is smaller than the one for Model M1 and indicates that for Model M2, the selected facilities appear to be closer in distance to the demand points than in Model M1.

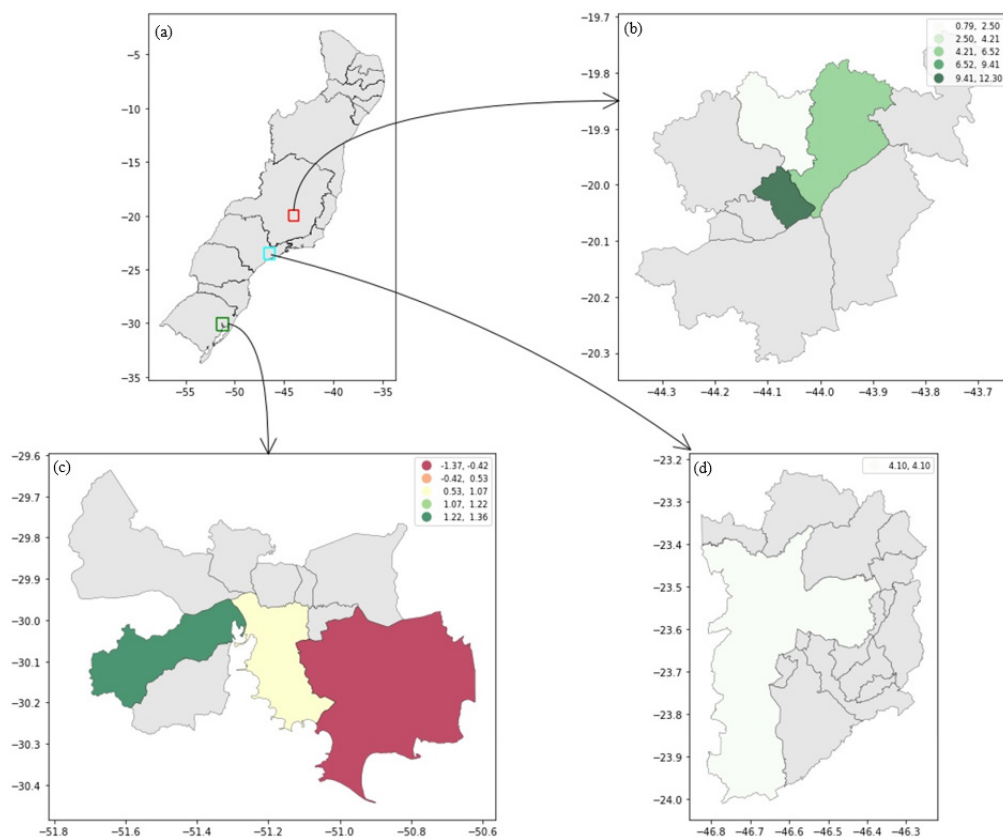
## 5.2. Evaluation of aggregation effect in instances of Set 2

The solutions for unaggregated and aggregated models of Set 2 instances are shown in Table 4. Regarding the objective function value, aggregated models are overestimating the unaggregated ones in instances 2-9 to 2-9000. Model M1 presents the largest gap between aggregated and unaggregated models in Set 2 solutions, equivalent to 0.41%. The objective function value rises as the fixed cost of opening new facilities increases.

The solutions of instances 2-9 and 2-4500 are the ones that present fewer open facilities in the optimal solution than in the unaggregated model. Again, the decision to open fewer facilities requires careful consideration of the corresponding elevation in transportation costs. This is underscored by the fact that the objective function value for aggregated models has shown an increase in such instances. For example, for instances 2-90 and 2-9000, the differences in the objective function values indicate that the set of open facilities in the aggregated models is different from the one in the unaggregated. Different sets of open facilities due to aggregation are responsible for the gaps in the objective function values. Instances that yield equal objective function values in both aggregated and unaggregated models imply identical optimal solutions in terms of open facilities.



**Figure 1.** Transportation cost differences for Instance 1-9 with demand points aggregated for Model M1 in (b) Brazil highlighting the metropolitan areas of (a) Belo Horizonte, (c) Salvador, (d) Porto Alegre, (e) São Paulo, and (f) Rio de Janeiro



**Figure 2.** Transportation cost differences for Instance 1-9 with demand points aggregated for Model M2 in (a) Brazil highlighting the metropolitan areas of (b) Belo Horizonte, (c) Porto Alegre, and (d) São Paulo

Running Gurobi considering the unaggregated model required, on average, 163,432 simplex iterations to reach the optimal solution, while Models M1 and M2 needed only 6,683 and 7,106 iterations, respectively. Again, the unaggregated model required more iterations to find its optimal solution, which also impacted the total processing time. On average, the aggregated models took 1 second to find their optimal solutions, while the unaggregated model took 165 seconds (approximately 3 minutes) to determine the optimal solutions. This finding highlights the capacity of aggregation in finding optimal or near-optimal solutions by means of aggregating input demand data.

As with Set 1, to identify regions in space where there are differences between the results for the aggregated and unaggregated mathematical models, we apply the DBSCAN algorithm. The *Eps* parameter value is 0.017 and was calculated using the nearest neighbors' distance between demand points. DBSCAN identified that 18 municipalities have more than one cluster in their territory, 10 municipalities are in clusters with demand points from other municipalities, and 10 have more than one cluster and clusters with demand points from other municipalities. The remaining municipalities did not meet the criteria and were considered neither large municipalities nor belonging to a metropolitan area.

We selected instance 2-4500 of Set 2 to analyze because it is the one in which the results of both aggregated models exhibit the same gap with respect to the result of the unaggregated model. Thus, for each municipality, transportation costs with respect to unaggregated and aggregated models are calculated to be compared.

The location of demand points allocated to different facilities in Model M1 is depicted in Figure 3. It can be observed that there are demand points allocated to different facilities in five distinct regions. Differences in transportation costs range between -84.99% and 662.90%, which means that there are municipalities where the transportation costs of the aggregated models either underestimate or overestimate the actual costs in comparison with the unaggregated model.

In its turn, Figure 4 depicts the locations of demand points allocated to different facilities in the results of aggregated and unaggregated models and the differences in transportation costs for Model M2. It can be observed that there are demand points allocated to different facilities in only three regions. The differences in transportation costs range between -7.71% and 110.31%. The differences in transportation costs for Model M1 are greater than the ones in Model M2, which indicates that the open facilities are closer to the demand points in Model M2.

It is important to observe that Model M1 exhibits a notably smaller objective function value compared to Model M2, even while having one less open facility. This observation implies that within a pool of alternative solutions, the solution derived from Model M1 would be superior to that of Model M2. However, if we analyze the transportation costs range, the difference between aggregated models is substantial. This is due to the location of the facilities selected in the aggregated models. The positive values of the range indicated that the facility opened in the aggregated model is more distant to a set of demand points than the one in the unaggregated model. On the other hand, the negative values indicate the opposite: the open facility is closer to the set of demand points in the aggregated model than in the unaggregated. In our problem, it means that for Model M2 the open facilities are generally closer to the demand points than the ones in Model M1.

**Table 4:** Set 2 solutions

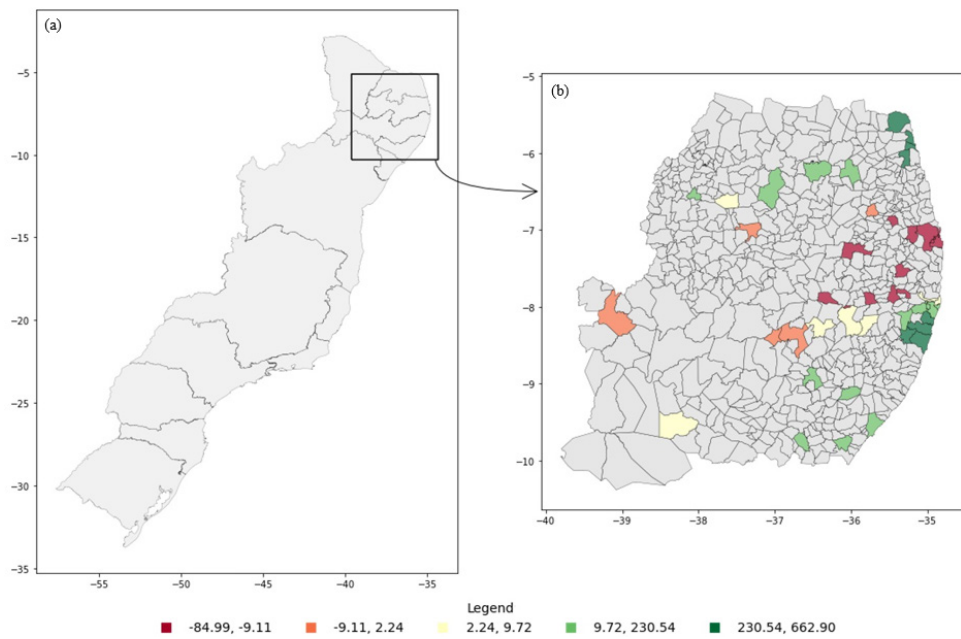
Instance/Model	2-9			2-90			2-900		
	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2
Objective Function	49,615.16	49,819.17	49,616.10	55,118.82	55,191.36	55,118.82	82,236.44	82,236.44	82,236.44
Difference	-	0.41%	<b>0.00%</b>	-	<b>0.00%</b>	0.00%	-	0.00%	0.00%
# Open Facilities	117	103	117	76	76	76	33	33	33
Processing time	21	2	1	51	1	1	163	1	1
Explored nodes	209	0	0	1,284	0	0	1,655	0	0
Simplex iterations	5,668	195	236	32,979	612	813	43,423	2,400	2,631

Instance/Model	2-1800			2-4500			2-9000		
	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2
Objective Function	99,888.29	99,888.29	99,888.29	137,400.95	137,413.37	137,417.77	182,234.75	182,663.03	182,234.75
Difference	-	0.00%	0.00%	-	0.01%	0.01%	-	0.24%	0.00%
# Open Facilities	22	22	22	15	14	15	10	10	10
Processing time	156	1	1	127	1	1	135	1	1
Explored nodes	1,828	0	0	1,707	0	0	1,758	0	0
Simplex iterations	59,445	3,395	3,671	95,226	5,197	5,520	120,268	7,124	7,291

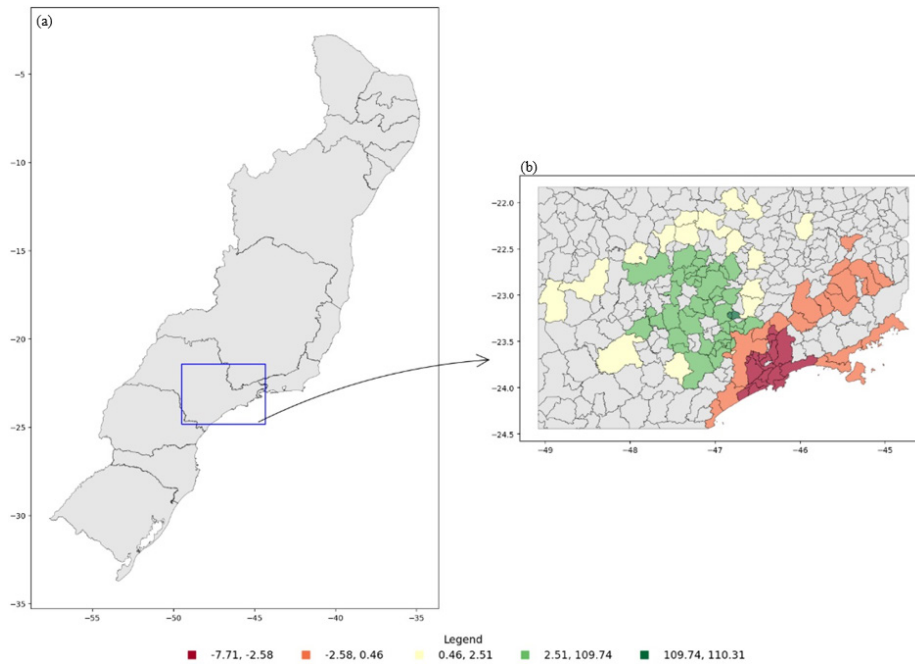
Instance/Model	2-18000			2-45000			2-90000		
	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2	Unaggregated	Model M1	Model M2
Objective Function	248,724.71	248,724.71	248,724.71	364,562.14	364,562.14	364,562.14	483,963.76	483,963.76	483,963.76
Difference	-	0.00%	0.00%	-	0.00%	0.00%	-	0.00%	0.00%
# Open Facilities	8	8	8	4	4	4	3	3	3
Processing time	419	1	2	261	2	2	150	2	2
Explored nodes	5,742	0	0	3,477	0	0	1,562	0	0
Simplex iterations	221,939	9,068	9,335	370,581	13,283	13,482	521,360	20,495	20,976



**Figure 3.** Transportation cost difference for Instance 2-4500 with demand points aggregated for Model M1 in (a) Brazil highlighting the (b) municipalities in Northeast Region.

The different set of open facilities is associated with the aggregation method and the location of the clusters' center of gravity. Through the analysis of both Set 1 and 2 instances, it was possible to observe that the aggregation method takes a central role in the decision to open facilities to supply the demand. Model M2 clustering points with the K-means algorithm provided overall objective function values closer to the unaggregated models, and open facilities closer to the demand points than Model M1. This supports our hypothesis that in a situation with thousands of demand points unevenly spread across a vast area, using an algorithm such as K-means, which does not take into consideration the municipality limits, can lead to better results as not only the set of selected facilities is generally the same as in

the unaggregated model, but also the allocation of demand to them is identical. On the other hand, model M1 yields more location errors, as different facilities are selected.



**Figure 4.** Transportation cost difference for Instance 2-4500 with demand points aggregated with Method M2 in (a) Brazil highlighting the (b) municipalities in Southeast Region.

One possibility to avoid this difference range is by improving the definition of candidates in the modeling phase. In the case of Model M2, it is possible to increase the number  $K$  of clusters when aggregating demand data. However, it may come into an aggregation paradox, as increasing the number of clusters does not necessarily lead to an improvement in the value of the objective (Erkut and Bozkaya, 1999).

## 6. CONCLUDING REMARKS

There is often a trade-off between the level of detail and the ease of solving a facility location mathematical model. Also, as pointed out by Goodchild (1979), while uncertainty can be present in facility location models in different ways, real-world problems may require the aggregation of the set of demand points for the sake of manageability, as they may encompass a large number. The idea behind the aggregation is to reduce the number of demand points to be small enough to allow obtaining the optimal solutions by means of an exact optimization model in reasonable computing times. However, such aggregation may reduce the accuracy of the model, as it may affect both the location of the selected facilities as well as the allocation of the demand points to them.

In this paper, we have investigated two demand point aggregation methods for the Fixed Charge Facility Location Problem (FCFL) to assess the impact on the solutions due to the use of aggregated data. Using real-world data, we generated 18 different instances of the FCFL based on two datasets (Set 1 and 2) that represent a real B2B distribution in Brazil to conduct the analysis.



For both sets, results indicate that demand point aggregation significantly reduced the processing time and computational effort. The results showed that for instances with higher fixed costs, which lead to a lower number of facilities that cover larger areas, aggregation effects are not significant since all models with aggregated demand points present the same value for the objective function cost as the unaggregated model.

In general, all differences in aggregated models are related to the clustering method employed, especially for instances in which the fixed cost of opening a facility is low in comparison with the transportation costs. For these instances, Model M2 allowed the identification of more open facilities, especially in areas with a higher number of demand points. Also, for Model M2, the differences in transportation costs range are smaller than for Model M1, indicating that open facilities are closer to demand points than in Model M1. Our results using real-world instances with thousands of demand points unevenly spread across a vast area show that employing a clustering algorithm such as K-means can lead to better results, oftentimes equal to those obtained by solving the respective unaggregated models, which require longer processing times. The results may also suggest that municipality boundaries, typically defined in a distant past, for some reason and with some purpose that is oftentimes irrelevant to a location study, do not yield high-quality solutions for the aggregation of demand points that are represented by the municipality centroid. On the other hand, model accuracy is barely affected when aggregation is performed by K-means, even in the case the number of clusters  $K$  is equal to the number of municipalities. It was also noted that for municipalities outside larger metropolitan areas, the model was less accurate, as demonstrated in the larger differences in instance 2-4500.

We highlight some research aspects that we believe to be worth pursuing. The comparison of the solutions of both clustering methods with heuristics for larger instances that could not be solved to optimality in reasonable times with the aim to assess in which cases each one of them would perform better is a relevant extension that can improve the discussion on whether more attention to the modeling should be paid, or if one should spend more time improving algorithms. This is particularly true in the context of deliveries to customers in a B2C distribution system in an urban context, such as in e-commerce, in which home addresses can exceed a million. However, it would require significant additional effort to be addressed properly which is beyond the scope of this study. Another suggestion is related to the usage of already aggregated data, such as census data since the representation can lead to imprecise results. Another promising extension would be to compare clustering methods for other related location problems, such as, for instance, the capacitated facility location problem. We leave them all as topics for future investigation.

#### ACKNOWLEDGEMENTS

The first and the second authors acknowledge Brazil's CNPq (National Council for Scientific and Technological Development) financial support [grants numbers 133773/2019-1 and 309424/2018-6, respectively].

#### REFERENCES

- Aardal, K. (1998) Capacitated facility location: separation algorithms and computational experience. *Mathematical Programming, Series B*, v. 81, n. 2, p. 149-175. DOI: 10.1007/BF01581103.
- Aardal, K.; Y. Pochet and L.A. Wolsey (1995) Capacitated facility location: valid inequalities and facets. *Mathematics of Operations Research*, v. 20, n. 3, p. 562-582. DOI: 10.1287/moor.20.3.562.

- Andersson, G.; R.L. Francis; T. Normark et al. (1998) Aggregation method experimentation for large-scale network location problems. *Location Science*, v. 6, n. 1-4, p. 25-39. DOI: 10.1016/S0966-8349(98)00045-X.
- Barahona, F. and F.A. Chudak (2005) Near-optimal solutions to large-scale facility location problems. *Discrete Optimization*, v. 2, n. 1, p. 35-50. DOI: 10.1016/j.disopt.2003.03.001.
- Barcelo, J.; Å. Hallefjord; E. Fernandez et al. (1990) Lagrangean relaxation and constraint generation procedures for capacitated plant location problems with single sourcing. *OR-Spektrum*, v. 12, n. 2, p. 79-88. DOI: 10.1007/BF01784983.
- Beasley, J. (1990) OR-Library: distributing test problems by electronic mail. *The Journal of the Operational Research Society*, v. 41, n. 11, p. 1069-1072. DOI: 10.1057/jors.1990.166.
- Brasil (2021) Resolução nº 5.949, de 13 de julho de 2021. Altera o Anexo II da Resolução ANTT nº 5.867, de 14 de janeiro de 2020, em razão do disposto nos parágrafos 1º e 2º do artigo 5º da Lei nº 13.703, de 8 de agosto de 2018. *Diário Oficial da República Federativa do Brasil*. Seção 1, Brasília.
- Bueno, A. (2019) Preço Médio do Aluguel de Salas e Conjuntos Comerciais Acelera em Maio. Available at: <<https://fipezap.zapimoveis.com.br/preco-medio-do-aluguel-de-salas-e-conjuntos-comerciais-acelera-em-maio/>> (accessed 10/16/2023).
- Cebecauer, M. and L. Buzna (2017) A versatile adaptive aggregation framework for spatially large discrete location-allocation problems. *Computers & Industrial Engineering*, v. 111, p. 364-380. DOI: 10.1016/j.cie.2017.07.022.
- Daskin, M.S. (2013) *Network and Discrete Location* (2nd ed.). Hoboken: John Wiley & Sons, Ltd.
- Daskin, M.S.; A.E. Haghani; M. Khanal et al. (1989) Aggregation effects in maximum covering models. *Annals of Operations Research*, v. 18, n. 1, p. 113-139. DOI: <http://dx.doi.org/10.1007/BF02097799>.
- De Armas, J.; A.A. Juan; J.M. Marquès et al. (2017) Solving the deterministic and stochastic uncapacitated facility location problem: from a heuristic to a simheuristic. *The Journal of the Operational Research Society*, v. 68, n. 10, p. 1161-1176. DOI: 10.1057/s41274-016-0155-6.
- Erkut, E. and B. Bozkaya (1999) Analysis of aggregation errors for the  $p$ -median problem. *Computers & Operations Research*, v. 26, n. 10-11, p. 1075-1096. DOI: 10.1016/S0305-0548(99)00021-0.
- Erlenkotter, D.A. (1978) Dual-based procedure for uncapacitated facility location. *Operations Research*, v. 26, n. 6, p. 992-1009. DOI: 10.1287/opre.26.6.992.
- Ester, M.; H.-P. Kriegel; J. Sander et al. (1996) A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD'96)* (Portland, OR). Washington, DC: AAAI Press, p. 226-231. DOI: 10.5555/3001460.3001507.
- Fischetti, M.; I. Ljubić and M. Sinnl (2016) Benders decomposition without separability: a computational study for capacitated facility location problems. *European Journal of Operational Research*, v. 253, n. 3, p. 557-569. DOI: 10.1016/j.ejor.2016.03.002.
- Francis, R.L.; T.J. Lowe; M.B. Rayco et al. (2009) Aggregation error for location models: survey and analysis. *Annals of Operations Research*, v. 167, n. 1, p. 171-208. DOI: 10.1007/s10479-008-0344-z.
- Galvão, R. and L. Raggi (1989) A method for solving to optimality uncapacitated location problems. *Annals of Operations Research*, v. 18, n. 1, p. 225-244. DOI: 10.1007/BF02097805.
- Gao, X. (2021) A location-driven approach for warehouse location problem. *The Journal of the Operational Research Society*, v. 72, n. 12, p. 2735-2754. DOI: 10.1080/01605682.2020.1811790.
- Goodchild, M.F. (1979) The aggregation problem in location-allocation. *Geographical Analysis*, v. 11, n. 3, p. 240-255. DOI: 10.1111/j.1538-4632.1979.tb00692.x.
- Gurobi (2022) Gurobi Optimizer Reference Manual. Available at: <<https://www.gurobi.com>> (accessed 10/16/2023).
- Hakli, H. and Z. Ortacay (2019) An improved scatter search algorithm for the uncapacitated facility location problem. *Computers & Industrial Engineering*, v. 135, p. 855-867. DOI: 10.1016/j.cie.2019.06.060.
- Han, J.; M. Kamber and J. Pei (2012) *Data Mining*. Boston: Morgan Kaufmann.
- Hansen, P.; J. Brimberg; D. Urošević et al. (2007) Primal-dual variable neighborhood search for the simple plant-location problem. *INFORMS Journal on Computing*, v. 19, n. 4, p. 552-564. DOI: 10.1287/ijoc.1060.0196.
- Hillsman, E.L. and R. Rhoda (1978) Errors in measuring distances from populations to service centers. *The Annals of Regional Science*, v. 12, n. 3, p. 74-88. DOI: 10.1007/BF01286124.
- Hodgson, M.J.; F. Shmulevitz and M. Körkel (1997) Aggregation error effects on the discrete-space  $p$ -median model: the case of Edmonton, Canada. *Canadian Geographer*, v. 41, n. 4, p. 415-428. DOI: 10.1111/j.1541-0064.1997.tb01324.x.
- Irawan, C.A. and S. Salhi (2015) Aggregation and non aggregation techniques for large facility location problems-a survey. *Yugoslav Journal of Operations Research*, v. 25, n. 3, p. 313-341. DOI: 10.2298/YJOR140909001I.
- Irawan, C.A.; S. Salhi; M. Luis et al. (2017) The continuous single source location problem with capacity and zone-dependent fixed cost: models and solution approaches. *European Journal of Operational Research*, v. 263, n. 1, p. 94-107. DOI: 10.1016/j.ejor.2017.04.004.

- Jacobs-Crisioni, C.; P. Rietveld and E. Koomen (2014) The impact of spatial aggregation on urban development analyses. *Applied Geography*, v. 47, p. 46-56. DOI: 10.1016/j.apgeog.2013.11.014.
- Janjevic, M.; M. Winkenbach and D. Merchán (2019) Integrating collection-and-delivery points in the strategic design of urban last-mile e-commerce distribution networks. *Transportation Research Part E: Logistics and Transportation Review*, v. 131, p. 37-67. DOI: 10.1016/j.tre.2019.09.001.
- Khumawala, B. (1972) An efficient branch and bound algorithm for the warehouse location problem. *Management Science*, v. 18, n. 12, p. B-635-B-749. DOI: 10.1287/mnsc.18.12.B718.
- Kratka, J.; D. Tošić; V. Filipović et al. (2001) Solving the simple plant location problem by genetic algorithm. *Operations Research*, v. 35, n. 1, p. 127-142. DOI: 10.1051/ro:2001107.
- Laporte, G.; S. Nickel and F. Saldanha-da-Gama (2019) *Location Science*. Cham: Springer. DOI: 10.1007/978-3-030-32177-2.
- Levin, Y. and A. Ben-Israel (2004) A heuristic method for large-scale multi-facility location problems. *Computers & Operations Research*, v. 31, n. 2, p. 257-272. DOI: 10.1016/S0305-0548(02)00191-0.
- MacQueen, J. (1967) Some methods for classification and analysis of multivariate observations. In Le Cam, L.M. and J. Neyman (eds.) *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability* (v. 1). Oakland: University of California Press, p. 281-297.
- Melo, M.; S. Nickel and F. Saldanha-da-Gama (2009) Facility location and supply chain management – a review. *European Journal of Operational Research*, v. 196, n. 2, p. 401-412. DOI: 10.1016/j.ejor.2008.05.007.
- O’Kelly, M.E. (1992) A clustering approach to the planar hub location problem. *Annals of Operations Research*, v. 40, n. 1, p. 339-353. DOI: 10.1007/BF02060486.
- Paul, J.A. and L. Macdonald (2016) Location and capacity allocations decisions to mitigate the impacts of unexpected disasters. *European Journal of Operational Research*, v. 251, n. 1, p. 252-263. DOI: 10.1016/j.ejor.2015.10.028.
- Rahmah, N. and I.S. Sitanggang (2016) Determination of optimal Epsilon (Eps) value on DBSCAN algorithm to clustering data on peatland hotspots in Sumatra. *IOP Conference Series: Earth and Environmental Science*, v. 31, n. 1, p. 012012. DOI: 10.1088/1755-1315/31/1/012012.
- ReVelle, C.S.; H.A. Eiselt and M.S. Daskin (2008) A bibliography for some fundamental problem categories in discrete location science. *European Journal of Operational Research*, v. 184, n. 3, p. 817-848. DOI: 10.1016/j.ejor.2006.12.044.
- Rodero, R. (2018) Gostaria de Saber Qual o Tamanho Mínimo de Uma Drogaria por Partes. Available at: <<https://guiadafarmacia.com.br/perguntas/qual-o-tamanho-minimo-de-uma-drogaria-por-partes/>> (accessed 10/16/2023).
- Saldanha-da-Gama, F. (2022) Facility location in logistics and transportation: an enduring relationship. *Transportation Research Part E: Logistics and Transportation Review*, v. 166, p. 102903. DOI: 10.1016/j.tre.2022.102903.
- Sankaran, J.K. (2007) On solving large instances of the capacitated facility location problem. *European Journal of Operational Research*, v. 178, n. 3, p. 663-676. DOI: 10.1016/j.ejor.2006.01.035.
- Tong, D. and R.L. Church (2012) Aggregation in continuous space coverage modeling. *International Journal of Geographical Information Science*, v. 26, n. 5, p. 795-816. DOI: 10.1080/13658816.2011.615748.
- Van Roy, T.A. (1986) Cross decomposition algorithm for capacitated facility location. *Operations Research*, v. 34, n. 1, p. 145-163. DOI: 10.1287/opre.34.1.145.
- Verter, V. (2011) Uncapacitated and capacitated facility location problems. In Eiselt, H.A. and V. Marianov (eds.) *Foundations of Location Analysis*. New York: Springer, p. 25-37. DOI: 10.1007/978-1-4419-7572-0\_2.
- Wolsey, L.A. (1998) *Integer Programming*. New York: John Wiley & Sons, Ltd.
- Zhao, P. and R. Batta (1999) Analysis of centroid aggregation for the Euclidean distance p-median problem. *European Journal of Operational Research*, v. 113, n. 1, p. 147-168. DOI: 10.1016/S0377-2217(98)00010-1.