



Article

A Method with Adaptive Graphs to Constrain Multi-View Subspace Clustering of Geospatial Big Data from Multiple Sources

Qiliang Liu , Weihua Huan and Min Deng *

Department of Geo-Informatics, Central South University, Changsha 410083, China

* Correspondence: dengmin@csu.edu.cn

Abstract: Clustering of multi-source geospatial big data provides opportunities to comprehensively describe urban structures. Most existing studies focus only on the clustering of a single type of geospatial big data, which leads to biased results. Although multi-view subspace clustering methods are advantageous for fusing multi-source geospatial big data, exploiting a robust shared subspace in high-dimensional, non-uniform, and noisy geospatial big data remains a challenge. Therefore, we developed a method with adaptive graphs to constrain multi-view subspace clustering of multi-source geospatial big data (agc2msc). First, for each type of data, high-dimensional and noisy original features were projected into a low-dimensional latent representation using autoencoder networks. Then, adaptive graph constraints were used to fuse the latent representations of multi-source data into a shared subspace representation, which preserved the neighboring relationships of data points. Finally, the shared subspace representation was used to obtain the clustering results by employing a spectral clustering algorithm. Experiments on four benchmark datasets showed that agc2msc outperformed nine state-of-the-art methods. agc2msc was applied to infer urban land use types in Beijing using the taxi GPS trajectory, bus smart card transaction, and points of interest datasets. The clustering results may provide useful calibration and reference for urban planning.

Keywords: multi-view subspace clustering; geospatial big data; shared nearest neighbor graph; social sensing



Citation: Liu, Q.; Huan, W.; Deng, M. A Method with Adaptive Graphs to Constrain Multi-View Subspace Clustering of Geospatial Big Data from Multiple Sources. *Remote Sens.* **2022**, *14*, 4394. <https://doi.org/10.3390/rs14174394>

Academic Editors: Qiming Zhou, Jianfeng Li, Meng Gao and Bin Chen

Received: 20 July 2022

Accepted: 1 September 2022

Published: 3 September 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Multi-source geospatial big data have become increasingly available in the current era of big data, such as taxi GPS trajectories [1], smart card transactions [2], mobile phone data [3], social media check-in records [4], and points of interests (POIs) [5]. Geospatial big data provides a new opportunity for understanding the “human-earth” relationship [6]. Clustering geospatial big data are vital for describing urban structures and understanding the organization of cities [7]. For example, remote sensing techniques have been widely used for uncovering urban land use information based on physical characteristics of ground components (e.g., spectral, shape, and texture) [8]; however, remote sensing techniques are hard to capture the socioeconomic attributes and human dynamics that are highly related to urban land use [3]. In contrast, clustering of human mobility data can help understand urban land use information from the perspective of social function which is an important complement of remote sensing [6]. Clustering of geospatial big data are also useful for identifying urban functional structures and human activity patterns, which are useful for human-centric urban planning [9–11]. For example, the actual functions of a region may be inconsistent with the original zoning scheme designed by urban planners [12]. Clusters discovered from geospatial big data can reveal the urban function zones naturally formulated according to human activities, which may provide useful calibration for urban planners [9]. Clusters discovered from social media check-in records are also useful for identifying emergency events in a city, which are helpful for maintaining public safety [4].

Although clustering of geospatial big data has received attention in recent years, most existing studies focus on a single type of geospatial big data [11,13]. Owing to the bias of each type of geospatial big data, the clustering results obtained from single-source geospatial big data cannot provide a comprehensive view of urban structures [14]. A few studies have used a weighted average strategy to fuse multi-source geospatial big data [15,16]. Multi-source geospatial big data usually reflect different or overlapping dimensions of human activities. Without considering the shared and complementary information among different types of geospatial big data, the weighted average strategy may introduce unpredictable errors [17]. Multi-view subspace clustering has the potential to fuse the underlying complementary information of multi-source geospatial big data [18,19]; however, high-dimensional, non-uniform, and noisy geospatial big data bring two challenges [20–23]: (1) the quality of the low-dimensional subspace is substantially influenced by the redundant features and noise in the original data; and (2) neighboring relationships of data points in high-dimensional and non-uniform original data space are difficult to preserve in a low-dimensional subspace. Therefore, existing multi-view subspace clustering methods are highly likely to generate an inaccurate subspace, which degrades the clustering performance. To overcome the above challenges, this study developed a method with adaptive graphs to constrain multi-view subspace clustering of geospatial big data from multiple sources (agc2msc). The main contributions of this work include the following three aspects:

- (1) We used a multi-view learning strategy to fuse the information embedded in multi-source geospatial big data. Compared with the weighted average strategy, multi-view subspace clustering is more suitable for integrating different and/or overlapping dimensions of human activities reflected by multi-source geospatial big data.
- (2) We used autoencoder networks [24] to map high-dimensional and noisy original geospatial big data into a latent representation. The latent representation of each type of geospatial big data was used to construct the low-dimensional subspace. Therefore, the influence of feature redundancy and noise on subspace construction can be reduced; moreover, the non-linear relationship between each type of data and its latent representation can be captured.
- (3) We used a shared nearest neighbor method [25] to construct adaptive graphs for high-dimensional, non-uniform, and noisy geospatial big data. The adaptive graphs can be used as constraints to obtain a more robust subspace shared by multi-source geospatial big data. Therefore, the quality of multi-view subspace clustering can be improved.

Experiments on four multi-view benchmark datasets showed that agc2msc outperformed nine state-of-the-art methods. A case study in Beijing showed that agc2msc is a powerful tool for inferring urban land use types from multi-source geospatial big data (i.e., taxi GPS trajectory, bus smart card transaction, and POI datasets). The clustering results may provide useful calibration and reference for urban planning.

2. Related Work

Most existing studies mainly focus on the clustering of a single type of geospatial big data, e.g., taxi GPS trajectories [1], social media check-in records [4], POIs [13], and mobile phone data [26]. After extracting clustering features from a certain type of geospatial big data, traditional clustering methods such as k-means [27], spectral clustering [28], and DBSCAN [29] are used to identify clusters. To consider the dynamic characteristic of geospatial big data, some online and incremental clustering methods are also currently available [4]; these methods are useful for understanding the organizations of cities from the perspective of social functions [7]. Despite these fruitful results, the bias of a single type of geospatial big data hinders the comprehensive understanding of urban structures [11,17]. To overcome this limitation, clustering of multi-source geospatial big data has received increasing attention in recent years. For example, some scholars [30] first combined the taxi trajectory data and public transit records to reveal human mobility patterns, then used POI features as prior knowledge to extract features of human mobility patterns, and finally

performed k-means on the extracted features. To consider the contributions of different types of geospatial big data, the weighted average strategy was employed to fuse the features of multi-source geospatial big data. The weights of different types of geospatial big data can be determined based on the proportions of total bus and cab ridership [15] or the entropy weight approach [16]. The weighted average methods can fuse the information of multi-source geospatial big data to a certain extent; however, they cannot incorporate complex interactions and correlations among multi-source geospatial big data. As shown in Figure 1, we can assume that the cone reflects the socioeconomic information that comprehensively describe the urban structures (i.e., the underlying structure of multi-source geospatial big data). In practice, this socioeconomic information is often embedded in different types of geospatial data (e.g., triangle and circle). Different types of geospatial data can be regarded as different views to observe socioeconomic information. The weighted average strategy does not capture the complementarity of multi-source geospatial big data. Therefore, the result of the weighted average strategy may be only a simple superposition of multiple features (i.e., the superposition of triangle and circle in Figure 1). Therefore, the underlying structure of multi-source geospatial big data cannot be reconstructed by using the weighted average strategy.

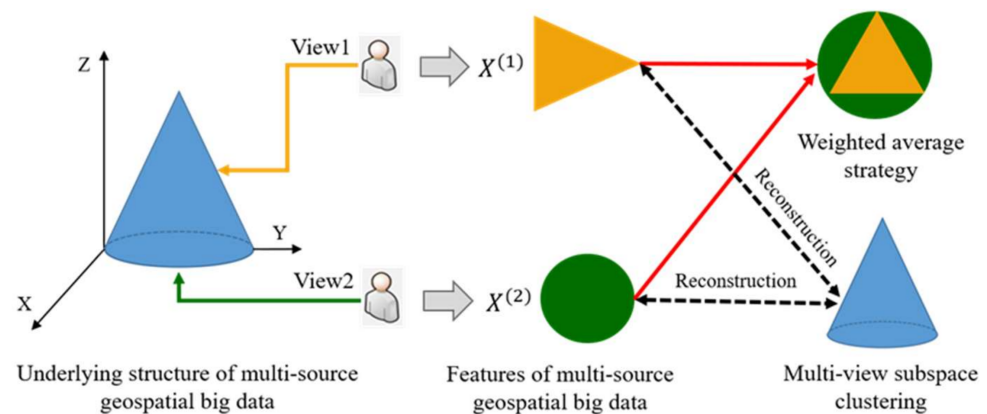


Figure 1. Illustration of the weighted average strategy and multi-view subspace clustering for fusing multi-source geospatial big data ($\{X^{(v)}\}_{v=1}^V$ ($V \geq 2$) refer to the features extracted from different types of geospatial big data).

Compared with the weighted average strategy, multi-view subspace clustering has the potential to reconstruct the underlying structure of multi-source geospatial data [18,19]. Multi-view subspace clustering assumes that multi-view data points are drawn from a shared low-dimensional subspace, rather than being uniformly distributed in the original space [31]. The features of each type of geospatial big data can be reconstructed from the shared subspace (the cone in the lower right corner of Figure 1). In theory, multi-view subspace clustering can fuse the shared and complementary information among different types of geospatial big data. Existing multi-view subspace clustering methods are mainly extensions of self-representation-based subspace clustering methods [32,33]. Self-representation-based subspace clustering assumes that each point x_i can be represented by a linear combination of other points x_j ($j \neq i$) [34–37]. Previous multi-view subspace clustering methods first calculate a subspace representation for each type of data and then combine the multiple subspace representations for clustering [21,31,38–40]. Although these methods can consider the shared and/or specific information of multi-source data, the subspaces reconstructed using the original data are not robust to redundant features and noise in the original data [41]. To address this limitation, latent multi-view subspace clustering methods have recently been developed [18,22]; these methods first use dimension reduction techniques to project the original data features into a latent representation, and then use the latent representation for subspace clustering. Although latent multi-view subspace clustering methods can boost the clustering performance of multi-source geospatial big

data, two challenges should be further addressed: (1) Existing method usually used a linear projection to transform the original data features into a latent representation [22,41,42]; however, the relationship between each type of data and its latent representation is usually non-linear [18,43]. Therefore, the inaccurate latent representations obtained by existing methods may degrade the clustering performance. (2) The neighboring relationships of data points in high-dimensional, non-uniform, and noisy original data are difficult to preserve in the shared subspace [36,44]. Some scholars have used neighbor graphs as constraints to preserve the neighboring relationships of data points in multi-view subspace clustering [21,45,46]; however, the neighbor graphs defined based on Euclidean distance and k-nearest neighbor cannot construct appropriate neighboring relationships for high-dimensional and non-uniform geospatial big data [47,48]. Therefore, existing methods are highly likely to generate an inaccurate subspace, which will reduce the clustering quality [49].

To overcome the above challenges, this study developed a method with adaptive graphs to constrain multi-view subspace clustering of geospatial big data from multiple sources.

3. Method

The framework of agc2msc is displayed in Figure 2. First, a data integration operation [50] should be performed to match multi-source geospatial big data to pre-defined spatial units. In this study, we matched each data point to traffic analysis zones according to their spatial location. Second, different types of features were extracted from multi-source geospatial big data, and a shared nearest neighbor graph was constructed to model the neighboring relationships of high-dimensional, non-uniform, and noisy geospatial big data. Third, an autoencoder network was utilized for encoding each type of original feature into a latent representation. Fourth, based on the self-representation property, the multiple latent representations were fused into a shared subspace representation under the constraint of the shared nearest neighbor graphs. The self-representation matrix and similarity matrix can be obtained. Finally, a spectral clustering algorithm [28] was used to obtain the clustering results with the similarity matrix. The main notations used in this paper are summarized in Table 1.

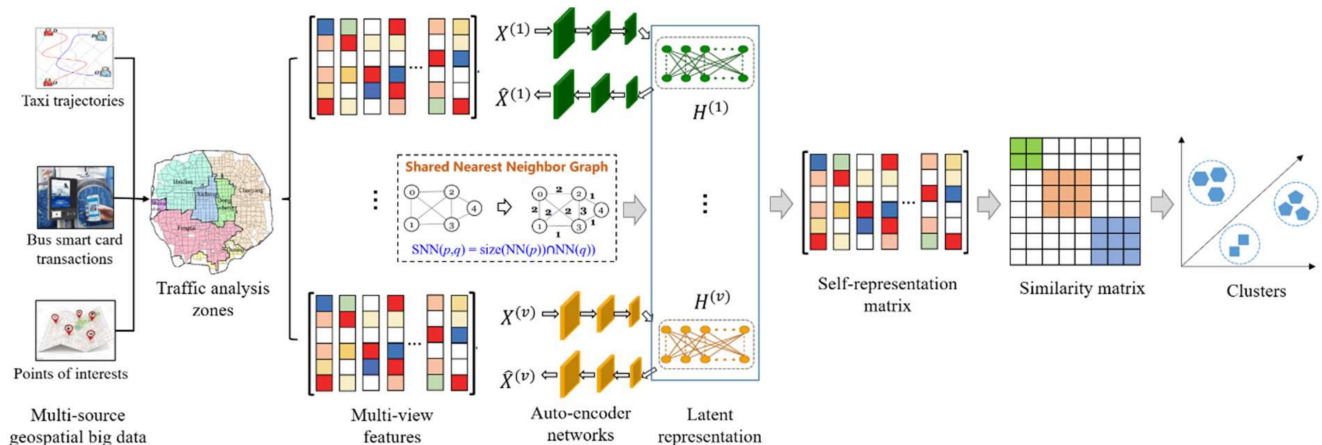


Figure 2. Framework of agc2msc ($\{\hat{X}^{(v)}\}_{v=1}^V$ ($V \geq 2$) refer to the reconstructed multi-view features; $\{H^{(v)}\}_{v=1}^V$ ($V \geq 2$) refer to the latent representations of multi-view features; $SNN(p, q)$ represents the shared nearest neighbor similarity between points p and q . $NN(p)$ and $NN(q)$ represent the k-nearest neighbors of p and q).

Table 1. Main notations and meanings through this paper.

| Notations | Meaning |
|---|--|
| O_w^i | The number of pick-ups in the i^{th} hour on weekdays |
| O_r^i | The number of pick-ups in the i^{th} hour on weekends |
| D_w^i | The number of drop-offs in the i^{th} hour on weekdays |
| D_r^i | The number of drop-offs in the i^{th} hour on weekends |
| $NN(x_i)$ | k nearest neighbors of x_i |
| N | The number of data points |
| V | The number of views |
| $W^{(v)}$ | The adjacency matrix of the v^{th} view |
| W | The unified adjacency matrix |
| w_{ij} | The shared nearest neighbor similarity between point x_i and point x_j |
| X | A multi-view dataset |
| $X^{(v)} \in \mathbb{R}^{d_v \times N}$ | Feature matrix of the v^{th} view |
| $\hat{X}^{(v)} \in \mathbb{R}^{d_v \times N}$ | Output of the decoder network of the v^{th} view |
| d_v | The dimension of the feature matrix in the v^{th} view |
| $H^{(v)}$ | Latent representation of the v^{th} view |
| $\theta_{(l,v)}$ | Combination of weights and bias in the m^{th} layer of the v^{th} view |
| M | The total number of layers in autoencoder networks |
| $f_i^{(m,v)}$ | Output of the m^{th} layer of the autoencoder in the v^{th} view |
| $W^{(m,v)}$ | The weight of the m^{th} layer of the autoencoder in the v^{th} view |
| $b^{(m,v)}$ | The bias of the m^{th} layer of the autoencoder in the v^{th} view |
| Z | Shared subspace representation matrix |
| z_i | The subspace representation of point x_i |
| D | Diagonal matrix |
| S | Similarity matrix |
| α, β | trade-off parameters |
| N_{TAZ} | The number of traffic analysis zones |
| I | Indicator function |
| FD | Frequency density |
| CR | Category rate |

3.1. Clustering Feature Extraction of Multi-Source Geospatial Big Data

Human mobility and POI data are two important types of geospatial big data [51]. For human mobility data, the temporal dynamics of boarding/de-boarding in each unit are frequently used to construct clustering features [1–3,52]. Several features have been developed for various applications. In this study, we inferred urban land use using multi-source geospatial big data. Existing work has found that the daily pick-up and drop-off combination vectors are the most suitable for revealing socioeconomic information on urban land use [1,19]. Therefore, this feature vector was used in this study.

$$[O_w^1, \dots, O_w^{16}, O_r^1, \dots, O_r^{16}, D_w^1, \dots, D_w^{16}, D_r^1, \dots, D_r^{16}], \quad (1)$$

where O_w^i and O_r^i denote the number of pick-ups in the i^{th} hour on weekdays and weekends, respectively, and D_w^i and D_r^i denote the number of drop-offs in the i^{th} hour on weekdays and weekends, respectively.

For POI data, deep-learning language models have been widely used to extract clustering features related to urban land use [5,13,30]. The entire study area is regarded as a corpus. Each spatial unit can be considered as a document. The contextual information of words can be regarded as geographical contextual information of POIs. The shortest path method, based on a greedy algorithm, was used to construct a spatial unit-based corpus [4]. In this study, to fully consider spatial heterogeneity, Word2vec [53] and Doc2vec [54] were used to embed POI data into clustering features. Word2vec was used to extract global geographic context information and Doc2vec was used to extract local geographic context

information; these two vectors were concatenated to construct a clustering feature (POI_{vec}) for each spatial unit.

$$POI_{vec} = POI_{Word2vec} \cup POI_{Doc2vec}, \quad (2)$$

where $POI_{Word2vec}$ represents the vector embedded by Word2vec and $POI_{Doc2vec}$ represents the vector embedded by Doc2vec.

3.2. Construction of Shared Nearest Neighbor Graph

To model the neighboring relationships for high-dimensional and non-uniform geospatial big data, we used the number of shared nearest neighbors to measure the similarity between data points.

$$SNN(x_i, x_j) = \text{size}(\text{NN}(x_i) \cap \text{NN}(x_j)), \quad (3)$$

where $\text{NN}(x_i)$ and $\text{NN}(x_j)$ represent the k nearest neighbors of x_i and x_j , respectively.

In high-dimensional space, direct similarity (e.g., Euclidean distance and cosine similarity) is not accurate because data in high dimensions are very sparse [47]. The shared nearest neighbor similarity is an indirect similarity that can effectively alleviate the problem of high dimensionality [55]. In addition, the shared nearest neighbor similarity can effectively alleviate the problem of varying densities. The shared nearest neighbor similarity can be adaptively adjusted according to the density of data points as it is dependent only on the number of neighbors that any two points share [56,57]. Figure 3 shows the neighbor graphs of a 2-dimensional dataset constructed using the distance threshold (Figure 3a), k -nearest neighbor strategy (Figure 3b), and shared nearest neighbor similarity (Figure 3c). The neighboring relationships of low-density points cannot be constructed using a distance threshold. Neighboring relationships constructed using the k -nearest neighbor strategy are usually inappropriate for outliers and points on the border of the clusters.

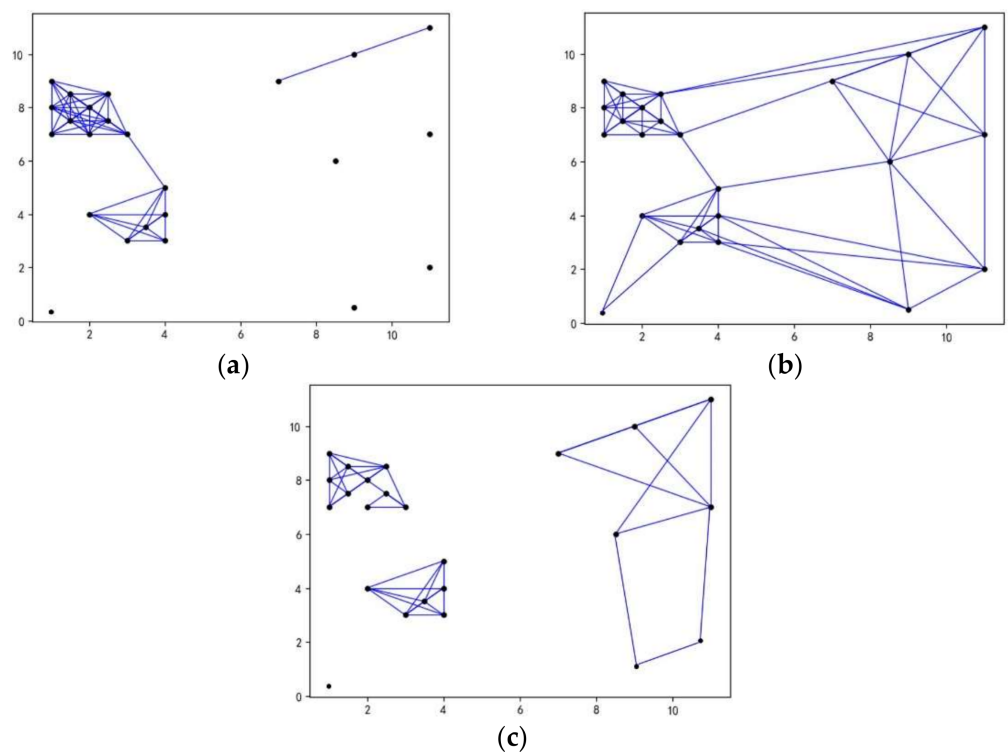


Figure 3. Illustration of the neighbor graph of data points. (a) Distance-based neighbor graph (threshold = 2.5). (b) k -nearest neighbor graph ($k = 5$). (c) Shared nearest neighbor graph ($k = 5$).

For each type of geospatial big data, a shared nearest neighbor graph was constructed to model the neighboring relationships of the data points. A unified graph was obtained by fusing multiple shared nearest neighbor graphs:

$$W = \frac{1}{V} \sum_{v=1}^V W^{(v)}, \quad (4)$$

where $W^{(v)}$ is the adjacency matrix of the v^{th} view and W is the unified adjacency matrix.

3.3. Latent Representation Based on Autoencoder Networks

In actual applications, the projection between an original feature and its latent representation is typically nonlinear. An autoencoder network [24] was used to obtain the latent representation for each type of clustering feature.

Given a multi-view dataset $X = \{X^{(1)}, X^{(2)}, \dots, X^{(V)}\}$, where $X^{(v)} \in \mathbb{R}^{d_v \times N}$ represents the feature matrix of the v^{th} view, d_v is the dimension of the v^{th} view, and N is the number of data points, the reconstruction loss of the autoencoder network can be formulated as:

$$L = \sum_v \left\| X^{(v)} - \hat{X}^{(v)} \right\|_F^2, \quad (5)$$

where $\hat{X}^{(v)}$ is the output of decoder network of the v^{th} view and $\|\cdot\|_F^2$ is the square of the Euclidean distance.

For each type of data, the output of the hidden layer of the autoencoder network is regarded as the latent representation, $H^{(v)} = [h_1^v, h_2^v, \dots, h_N^v]$:

$$h_n^v = g\left(\theta_{(M,v)} g\left(\theta_{(m,v)} \dots g\left(\theta_{(1,v)} [x_n^v; 1]\right)\right)\right), \quad (6)$$

where h_n^v is the latent representation of the n^{th} sample in the v^{th} view, $\theta_{(l,v)}$ is the combination of weights and bias of the m^{th} layer of the auto-coder networks in the v^{th} view, and $g(\cdot)$ is the activation function.

3.4. Multi-View Subspace Clustering with an Adaptive Graph Constraint

The shared and complementary information among different types of geospatial big data was fused based on the self-representation property [33]. Specifically, the latent representations of multi-source data were fused into a shared subspace representation under the constraint of the shared nearest neighbor graph:

$$\min_Z \sum_v L(H^{(v)}, H^{(v)}Z) + \Omega(Z) = \|H^{(v)} - H^{(v)}Z\|_F^2 + \Omega(Z), \quad (7)$$

where $L(\cdot)$ denotes the loss function, Z represents the shared subspace representation matrix, and $\Omega(\cdot)$ represents the graph regularization term that preserves the neighboring relationships of the data points in the shared subspace.

If a pair of points x_i and x_j are close in the original space, they should have similar subspace representations z_i and z_j , which can be formulated as the form $\|x_i - x_j\|_2^2 \rightarrow 0 \implies \|z_i - z_j\|_2^2 \rightarrow 0, \forall i \neq j$. Therefore, the regularization term $\Omega(Z)$ can be expressed as:

$$\Omega(Z) = \frac{1}{2} \sum_{i,j} w_{ij} \|z_i - z_j\|_2^2 = \text{tr}(Z^T D Z) - \text{tr}(Z^T W Z) = \text{tr}(Z^T L Z), \quad (8)$$

where w_{ij} denotes the shared nearest neighbor similarity between points x_i and x_j , $\|z_i - z_j\|_2^2$ represents the square of distances of z_i and z_j , and $\text{tr}(\cdot)$ represents the trace of the matrix. For a matrix $A \in \mathbb{R}^{n \times n}$, $\text{tr}(A) = \sum_{i=1}^n A_{ii}$. $L = W - D$ is a Laplacian matrix. W is the unified adjacency matrix, D is the diagonal matrix, and $d_{ii} = \sum_j w_{ij}$.

The objective function of agc2msc can be constructed by combining Equations (3)–(8):

$$\min_{\hat{X}^{(v)}, H^{(v)}, Z} \frac{1}{2} \sum_{v=1}^V \left(\|X^{(k)} - \hat{X}^{(v)}\|_F^2 + \alpha \|H^{(v)} - H^{(v)}Z\|_F^2 \right) + \beta \text{Tr}(Z^T LZ), \quad (9)$$

where α and β are trade-off parameters. The solution and optimization of Equation (8), are based on the Adam algorithm [58]. The details are presented in Appendix A. Z was used to construct the similarity matrix $S = \frac{1}{2}(|Z| + |Z|^T)$. A spectral clustering method was employed to obtain clusters using S . The silhouette coefficient was used to determine the optimal cluster number [59]. The pseudo code for the proposed method (Algorithm 1) is as follows:

Algorithm 1 The agc2msc method

Input:

multi-view dataset $X = \{X^{(1)}, X^{(2)}, \dots, X^{(V)}\}$, unified adjacency matrix W , parameters α and β .

Initial:

Learning rate: $lr = 0.001$

Optimizer: Adam

Epoch = 20,000

1: **While** pre-training not converged **do**:

2: Update $W^{(m,v)}$, $b^{(m,v)}$ and Z by formula (A3)–(A5) in Appendix A.

3: Obtain Z .

4: **End** pre-training.

5: **While** training not converged **do**:

6: Update $W^{(m,v)}$, $b^{(m,v)}$ and Z by formula (A3)–(A5) in Appendix A.

7: Obtain Z .

8: **End** training.

9: **Return** the shared subspace representation matrix Z .

Perform spectral clustering by employing the similarity matrix $S = \frac{1}{2}(|Z| + |Z|^T)$.

Output:

Clustering results.

4. Experiments

4.1. Benchmark Datasets

The performance of agc2msc was first evaluated using four multi-view clustering benchmark datasets: i.e., ORL, Yale, MSRCV1, and Caltech101-7 [23,39]. agc2msc was compared to nine state-of-the-art clustering methods. The quality of the clustering results was evaluated using the six indices, e.g., normalized mutual information (MNI), accuracy (ACC), F-score, Adjusted Rand index (AR), precision and recall. The results show that agc2msc outperforms the comparative methods on four benchmark datasets. The parameter sensitivity and model convergence of agc2msc were also evaluated. The experimental results are presented in Appendix B.

4.2. Case Study of Beijing Multi-Source Geospatial Big Data

4.2.1. Study Area and Dataset

agc2msc was applied to infer urban land-use types in Beijing from the perspective of social functions. The clustering results can reveal urban land use naturally formulated according to human activities, which may provide important complement of remote sensing [30]. Three types of geospatial big data were used in this study: taxi GPS trajectory, bus smart-card transactions, and POI datasets. The study area is located within the Fifth Ring Road of Beijing (Figure 4).

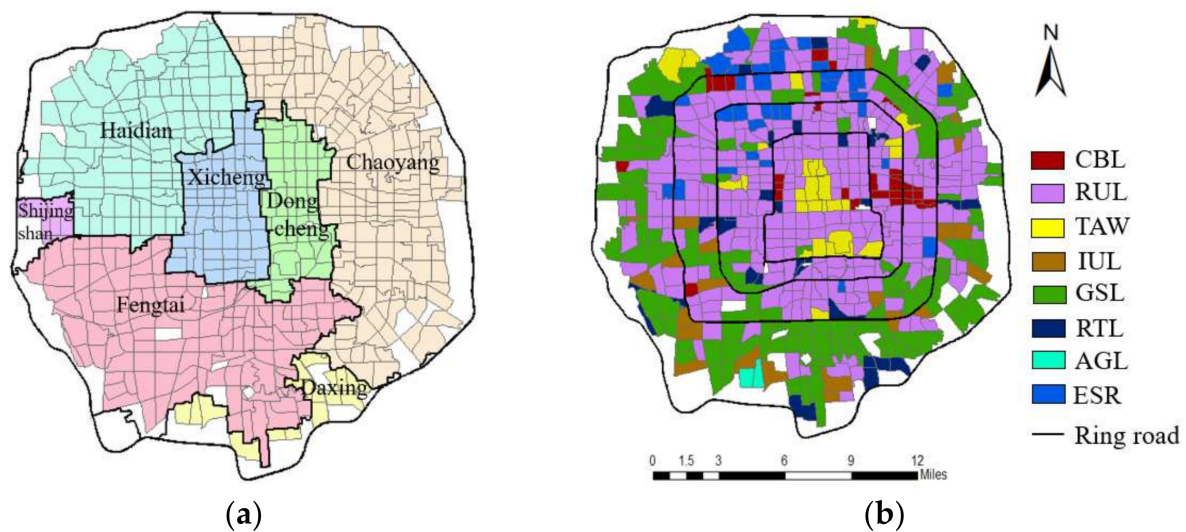


Figure 4. Study area. (a) Traffic analysis zones and administrative boundaries; (b) Beijing land use planning map (2017).

- (1) Traffic analysis zones. The study area was divided into 577 traffic analysis zones (Figure 4a). A traffic analysis zone is usually a socio-economically homogenous region that consists of one or more census blocks, block groups, or census tracts [60]. Existing work have found that traffic analysis zones are the suitable spatial units used in transportation and urban planning models [5,14,61]. Therefore, this study performs the clustering methods at the scale (or resolution) of traffic analysis zones. The traffic analysis zones were defined by the transport and urban planning authority, which were provided by Beijing Municipal Commission of Planning and Natural Resources.
- (2) Land use planning map: Figure 4b displays the governmental land-use map obtained from the Beijing Municipal Commission of Planning and Natural Resources. The current land classification (GB/T 21010-2017) identifies eight land use types: commercial and business land (CBL), residential land (RUL), tourist attraction and water (TAW), industrial land (IUL), green space land (GSL), road and transportation land (RTL), agricultural land (AGL), and education and scientific research land (ESR).
- (3) Taxi trajectory and bus smart-card transaction data: We collected GPS trajectories from more than 33,000 taxis and bus smart-card transaction data from 834 lines during the week (9:00–24:00 from 9 May 2016 to 15 May 2016). Each taxi trajectory contains the records of taxi ID, location, the status (occupied or not), and sampling time. We extracted the origin and destination points from each taxi trajectory. The numbers of taxi origin and destination pairs on workdays and weekends are 792,497 and 237,441, respectively. Each record of bus smart-card transaction data contains bus ID, the transaction time, pick-up station and drop-off station. For each bus smart-card transaction, the pick-up and drop-off stations were identified as the origin and destination points of that transaction. There are 14,157,913 and 4,157,948 bus origin and destination pairs on workdays and weekends, respectively. The origin and destination points of taxi trajectories and bus smart-card transactions were matched to traffic analysis zones according to their locations. The feature vectors constructed for the taxi GPS trajectory and bus smart-card transaction data had 64 dimensions;
- (4) POI data: POI data were collected from the 2017 Gaode Map. A total of 1,210,197 records were classified into 23 categories. Each POI record contained five essential attributes: name, ID, longitude, latitude, and category. For POI data, the information related to urban land use was extracted using two deep-learning language models, i.e., Word2vec and Doc2vec. We also matched POIs to traffic analysis zones according to the locations of POIs. A 64-dimensional feature vector was constructed for each traffic analysis zone.

4.2.2. Baseline Methods

We compared our method with the single-view spectral clustering [1], weighted average spectral clustering [15], and latent multi-view subspace clustering (gLMSC) [18] methods. The clustering features of the three comparative methods are the same as those of the proposed method. For the weighted average spectral clustering method, three similarity matrices, S_{taxi} , S_{bus} , S_{POI} , corresponding to the three types of data, were integrated using $S = \alpha_1 S_{\text{taxi}} + \alpha_2 S_{\text{bus}} + \alpha_3 S_{\text{POI}}$. In the experiment, α_1 , α_2 and α_3 were set to 3.66%, 96.34%, and 100%, respectively. We classified multi-source geospatial big data into two categories, i.e., human mobility data (taxi trajectories and bus smart-card transactions) and POI data. We set equal weights for human mobility data and POI data. Therefore, $\alpha_1 + \alpha_2 = \alpha_3 = 100\%$. For taxi trajectory data and bus smart-card transaction data, α_1 and α_2 were determined according to the proportion of taxi and bus ridership [15].

4.2.3. Clustering Results of agc2msc

For agc2msc, the silhouette coefficient reached its maximum value when the number of clusters was 10 (Figure 5a). To compare the clustering results, the clustering number of the other three methods was also set to 10. The clustering results of agc2msc are presented in Figure 5b. Two strategies were used to annotate the land-use types of identified clusters:

- (i) Frequency density (FD) and category rate (CR) of POIs in each cluster (Table 2):

$$FD_{ij} = \frac{\text{number of the } i_{\text{th}} \text{ category of POI in cluster } j}{\text{the area of cluster } j}, \quad (10)$$

$$CR_{ij} = \frac{\text{number of the } i_{\text{th}} \text{ category of POI in cluster } j}{\text{the number of POIs in cluster } j} \times 100\%, \quad (11)$$

- (ii) Arriving/leaving transition matrices: As shown in Figure 6, the horizontal axes represent the time over the day from 8:00 to 24:00, and the vertical axes represent the clusters for which passengers either arrive or leave. The colour for a grid represents the number of pick-ups or drop-offs in a cluster.

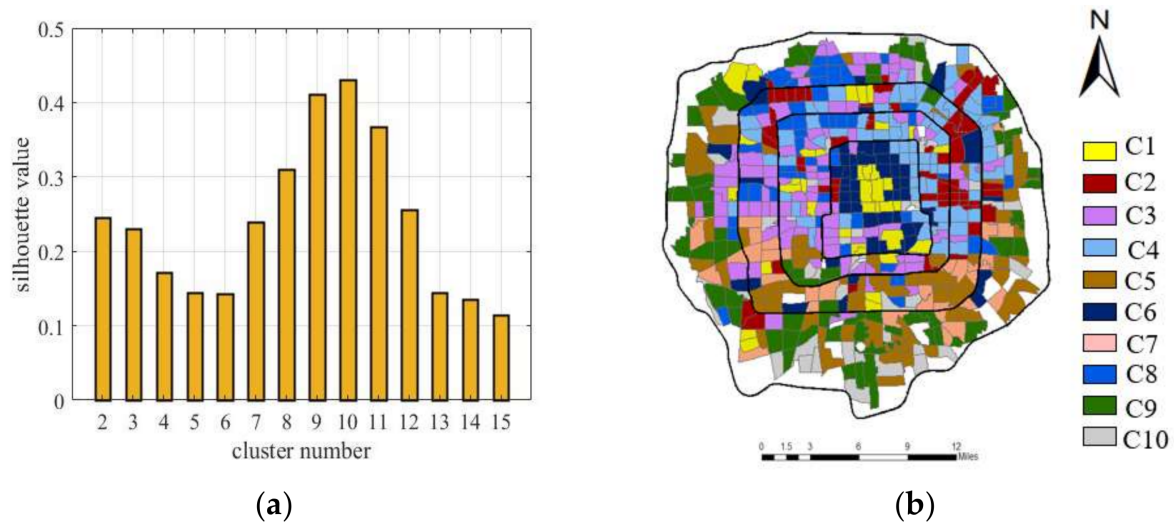


Figure 5. Clustering results. (a) Silhouette coefficient values with different cluster numbers; (b) Clustering results of the proposed method (ten clusters).

Table 2. Frequency Density (FD) and Category Rate (CR) of each cluster.

| | C1 | | C2 | | C3 | | C4 | | C5 | | C6 | | C7 | | C8 | | C9 | | C10 | |
|-------------------------|-------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|-------|--------|-------|--------|-------|--------|-------|--------|
| POI | FD | CR | FD | CR | FD | CR | FD | CR | FD | CR | FD | CR | FD | CR | FD | CR | FD | CR | FD | CR |
| Tourist attraction | 7.40 | 2.36% | 0.98 | 0.15% | 1.52 | 0.18% | 1.35 | 0.16% | 0.48 | 0.07% | 0.71 | 0.24% | 0.02 | 0.00% | 0.20 | 0.15% | 0.23 | 0.06% | 0.09 | 0.01% |
| Scenic spots | 13.73 | 4.39% | 1.63 | 0.24% | 1.88 | 0.22% | 2.42 | 0.28% | 1.23 | 0.17% | 1.16 | 0.40% | 0.95 | 0.20% | 0.49 | 0.38% | 5.41 | 1.30% | 0.22 | 0.03% |
| Hot place name | 0.19 | 0.06% | 0.01 | 0.00% | 0.02 | 0.00% | 0.01 | 0.00% | 0.04 | 0.01% | 0.00 | 0.00% | 0.04 | 0.01% | 0.00 | 0.00% | 0.01 | 0.00% | 0.01 | 0.00% |
| Cultural relics | 0.09 | 0.03% | 0.01 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.01 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% |
| Company/enterprise | 27.94 | 8.92% | 81.01 | 12.14% | 125.66 | 14.97% | 80.48 | 9.35% | 26.20 | 3.64% | 156.95 | 53.62% | 49.82 | 10.46% | 18.48 | 14.17% | 7.65 | 1.85% | 62.20 | 7.54% |
| Building | 1.65 | 0.53% | 4.20 | 0.63% | 5.75 | 0.69% | 5.79 | 0.67% | 1.70 | 0.24% | 7.17 | 2.45% | 1.67 | 0.35% | 0.36 | 0.28% | 2.10 | 0.51% | 0.60 | 0.07% |
| Shopping | 59.43 | 18.98% | 123.68 | 18.54% | 205.46 | 24.48% | 137.73 | 15.99% | 132.65 | 18.42% | 138.59 | 47.35% | 44.75 | 9.39% | 20.11 | 15.42% | 70.44 | 16.99% | 92.94 | 11.27% |
| Shopping mall | 0.97 | 0.31% | 7.67 | 1.15% | 6.39 | 0.76% | 4.34 | 0.50% | 2.98 | 0.41% | 2.20 | 0.75% | 0.99 | 0.21% | 0.13 | 0.10% | 1.75 | 0.42% | 0.40 | 0.05% |
| Theater | 0.25 | 0.08% | 1.32 | 0.20% | 1.82 | 0.22% | 0.97 | 0.11% | 0.82 | 0.11% | 1.72 | 0.59% | 0.36 | 0.08% | 0.07 | 0.05% | 2.32 | 0.56% | 0.60 | 0.07% |
| Accommodation | 6.93 | 2.21% | 7.63 | 1.14% | 36.71 | 4.37% | 23.21 | 2.69% | 16.73 | 2.32% | 20.95 | 7.16% | 8.52 | 1.79% | 2.40 | 1.84% | 12.96 | 3.13% | 2.10 | 0.25% |
| Catering service | 87.58 | 16.74% | 186.14 | 22.18% | 87.27 | 18.31% | 44.23 | 15.11% | 23.97 | 18.39% | 142.00 | 21.29% | 91.38 | 18.90% | 67.64 | 21.60% | 96.67 | 23.32% | 75.39 | 22.57% |
| Hotel | 21.43 | 2.49% | 28.19 | 3.36% | 12.53 | 2.63% | 6.42 | 2.19% | 3.45 | 2.65% | 20.26 | 3.04% | 13.77 | 1.46% | 9.24 | 2.95% | 15.76 | 3.80% | 16.60 | 2.12% |
| Dwelling | 14.67 | 4.68% | 10.69 | 1.60% | 51.81 | 6.17% | 35.75 | 4.15% | 34.62 | 4.81% | 41.76 | 14.27% | 3.70 | 0.78% | 4.88 | 3.74% | 20.06 | 4.84% | 4.15 | 0.50% |
| Courier service | 1.37 | 0.44% | 2.49 | 0.37% | 6.01 | 0.72% | 4.50 | 0.52% | 2.24 | 0.31% | 3.95 | 1.35% | 2.84 | 0.60% | 1.35 | 1.03% | 0.88 | 0.21% | 2.88 | 0.35% |
| Living service | 57.59 | 18.39% | 62.10 | 9.31% | 290.73 | 34.64% | 135.26 | 15.71% | 126.32 | 17.54% | 141.91 | 48.49% | 63.49 | 13.32% | 17.95 | 13.77% | 59.67 | 14.40% | 13.87 | 1.68% |
| Hair dressing | 15.98 | 5.10% | 29.12 | 4.37% | 80.77 | 9.62% | 29.77 | 3.46% | 34.72 | 4.82% | 37.33 | 12.75% | 18.77 | 3.94% | 3.27 | 2.51% | 13.86 | 3.34% | 13.97 | 1.69% |
| Health care treatment | 8.15 | 2.60% | 7.96 | 1.19% | 29.38 | 3.50% | 16.36 | 1.90% | 15.60 | 2.17% | 19.17 | 6.55% | 10.67 | 2.24% | 2.78 | 2.13% | 9.84 | 2.37% | 17.10 | 2.07% |
| Bank | 2.62 | 0.84% | 6.22 | 0.93% | 14.53 | 1.73% | 11.83 | 1.37% | 7.47 | 1.04% | 8.15 | 2.79% | 3.16 | 0.66% | 0.74 | 0.57% | 4.02 | 0.97% | 2.15 | 0.26% |
| Courier service | 5.09 | 1.62% | 7.42 | 1.11% | 14.95 | 1.78% | 8.75 | 1.02% | 8.15 | 1.13% | 3.44 | 1.17% | 6.35 | 1.33% | 2.14 | 1.64% | 5.60 | 1.35% | 4.74 | 0.57% |
| Moving company | 1.09 | 0.35% | 0.75 | 0.11% | 2.24 | 0.27% | 3.71 | 0.43% | 1.18 | 0.16% | 2.14 | 0.73% | 1.01 | 0.21% | 0.36 | 0.28% | 1.04 | 0.25% | 0.22 | 0.03% |
| Intermediary agency | 3.25 | 1.04% | 3.63 | 0.54% | 5.34 | 0.64% | 9.19 | 1.07% | 2.60 | 0.36% | 1.52 | 0.52% | 0.36 | 0.08% | 0.29 | 0.22% | 2.04 | 0.49% | 0.40 | 0.05% |
| Doorplate | 9.86 | 3.15% | 6.36 | 0.95% | 10.58 | 1.26% | 10.97 | 1.27% | 12.11 | 1.68% | 9.27 | 3.17% | 8.77 | 1.84% | 4.02 | 3.08% | 5.97 | 1.44% | 8.59 | 1.04% |
| Recreation place | 2.78 | 0.89% | 4.75 | 0.71% | 6.66 | 0.79% | 6.38 | 0.74% | 8.69 | 1.21% | 1.52 | 0.52% | 3.16 | 0.66% | 0.94 | 0.72% | 3.22 | 0.78% | 11.28 | 1.37% |
| Clothing factory | 11.30 | 3.61% | 20.14 | 3.02% | 68.46 | 8.16% | 28.60 | 3.32% | 28.80 | 4.00% | 4.97 | 1.70% | 42.39 | 8.90% | 1.14 | 0.87% | 15.03 | 3.63% | 5.18 | 0.63% |
| Industry | 1.25 | 0.40% | 0.59 | 0.09% | 0.67 | 0.08% | 0.33 | 0.04% | 0.33 | 0.05% | 1.05 | 0.36% | 1.84 | 0.39% | 0.45 | 0.35% | 0.39 | 0.09% | 0.33 | 0.04% |
| Educational service | 12.30 | 3.93% | 30.58 | 4.58% | 43.57 | 5.19% | 43.12 | 5.01% | 31.33 | 4.35% | 7.38 | 2.52% | 15.07 | 3.16% | 21.47 | 16.46% | 22.66 | 5.47% | 8.12 | 0.98% |
| Scientific institution | 4.28 | 1.37% | 3.27 | 0.49% | 8.54 | 1.02% | 7.16 | 0.83% | 4.61 | 0.64% | 2.19 | 0.75% | 1.29 | 0.27% | 4.63 | 3.55% | 1.81 | 0.44% | 1.17 | 0.14% |
| Sports leisure service | 9.24 | 2.95% | 20.26 | 3.04% | 28.19 | 3.36% | 23.77 | 2.76% | 26.60 | 3.69% | 6.42 | 2.19% | 12.53 | 2.63% | 3.45 | 2.65% | 15.76 | 3.80% | 3.32 | 0.40% |
| Natural place name | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.00 | 0.00% | 0.10 | 0.02% | 0.00 | 0.00% |
| Road ancillary facility | 2.12 | 0.68% | 2.76 | 0.41% | 3.94 | 0.47% | 4.10 | 0.48% | 3.72 | 0.52% | 1.34 | 0.46% | 1.73 | 0.36% | 1.05 | 0.80% | 5.25 | 1.27% | 11.13 | 1.35% |
| Sinopec | 0.19 | 0.06% | 0.06 | 0.01% | 0.02 | 0.00% | 0.07 | 0.01% | 0.06 | 0.01% | 0.26 | 0.09% | 0.02 | 0.00% | 0.05 | 0.04% | 0.17 | 0.04% | 0.20 | 0.02% |
| Gas station | 0.28 | 0.09% | 0.34 | 0.05% | 0.05 | 0.01% | 0.31 | 0.04% | 0.37 | 0.05% | 0.60 | 0.21% | 0.35 | 0.07% | 0.25 | 0.19% | 0.28 | 0.07% | 0.58 | 0.07% |
| Long-distance bus | 0.00 | 0.00% | 0.01 | 0.00% | 0.03 | 0.00% | 0.05 | 0.01% | 0.03 | 0.00% | 0.02 | 0.01% | 0.05 | 0.01% | 0.01 | 0.01% | 0.04 | 0.01% | 0.06 | 0.01% |
| Railway station | 0.00 | 0.00% | 0.04 | 0.01% | 0.05 | 0.01% | 0.01 | 0.00% | 0.06 | 0.01% | 0.05 | 0.02% | 0.02 | 0.00% | 0.02 | 0.01% | 0.10 | 0.02% | 0.20 | 0.02% |

Notes: If the FD value of a certain type of POI in cluster C_i is high, it indicates that the density of that kind of POI in C_i is high. If the CR value of a certain type of POI in cluster C_i is high, it indicates that the proportion of that kind of POI in all POIs in C_i is high.

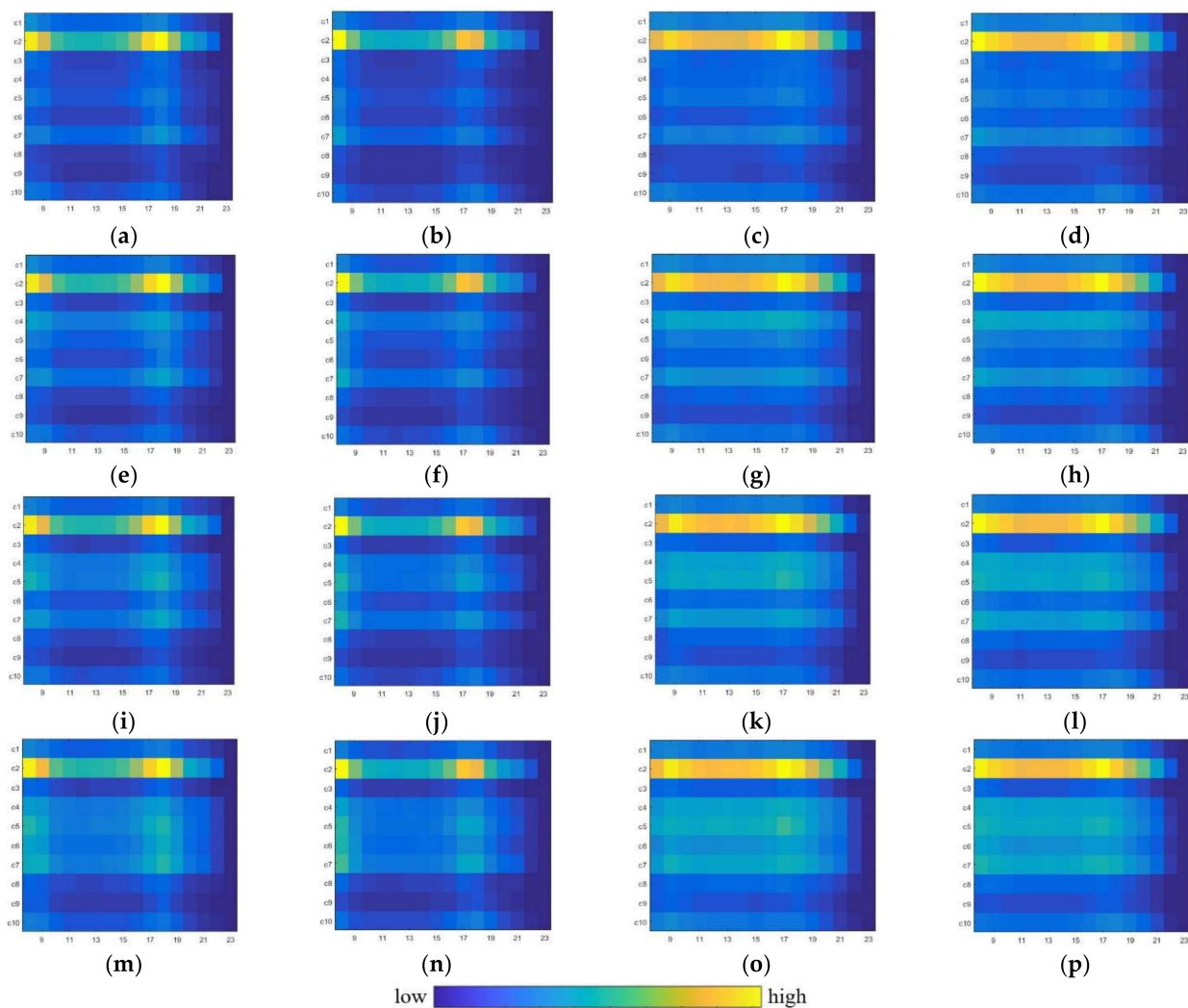


Figure 6. Arriving/leaving transition matrices. (a) arriving, C3, weekday; (b) leaving, c3, weekday; (c) arriving, c3, weekend; (d) leaving, c3, weekend; (e) arriving, c4, weekday; (f) leaving, c4, weekday; (g) arriving, c4, weekend; (h) leaving, c4, weekend; (i) arriving, c5, weekday; (j) leaving, c5, weekday; (k) arriving, c5, weekend; (l) leaving, c5, weekend; (m) arriving, c6, weekday; (n) leaving, c6, weekday; (o) arriving, c6, weekend; (p) leaving, c6, weekend. Horizontal axis: the time over the day from 8:00 to 24:00; Vertical axis: the clusters for which passengers either arrive or leave; Color of a grid: the number of pick-ups or drop-offs in a cluster.

Tourist Attraction and Water Areas (C1)

C1 was annotated as a tourist attraction and water area as it contains the largest number of tourists and attractions, scenic spots, hot place names, cultural relics, and aquariums among all clusters (Table 2). Figure 7a visualizes the intensity of several representative types of POIs such as toponymic address information, historical sites, scenic spots, cultural relics, and aquariums; it was found that most POIs in Beijing are centered at this cluster, such as Tiananmen Square, the Imperial Palace, Beijing Museum, Summer Palace, Old Summer Palace, and Temple of Heaven Park.

Commercial and business areas (C2)

C2 was annotated as a commercial and business area as it has an adequate POI configuration for companies, enterprises, buildings, shopping, catering services, theaters, hotels, and shopping malls (Table 2); moreover, Figure 7b visualizes several types of representative POIs in commercial areas, such as business trade, corporate enterprises,

restaurants, and bars; it was found that the prosperous core business circles are located in this cluster, such as Asian Games Village, Sanlitun, Zhongguancun, Wangjing, Central Business District, and Xidan.

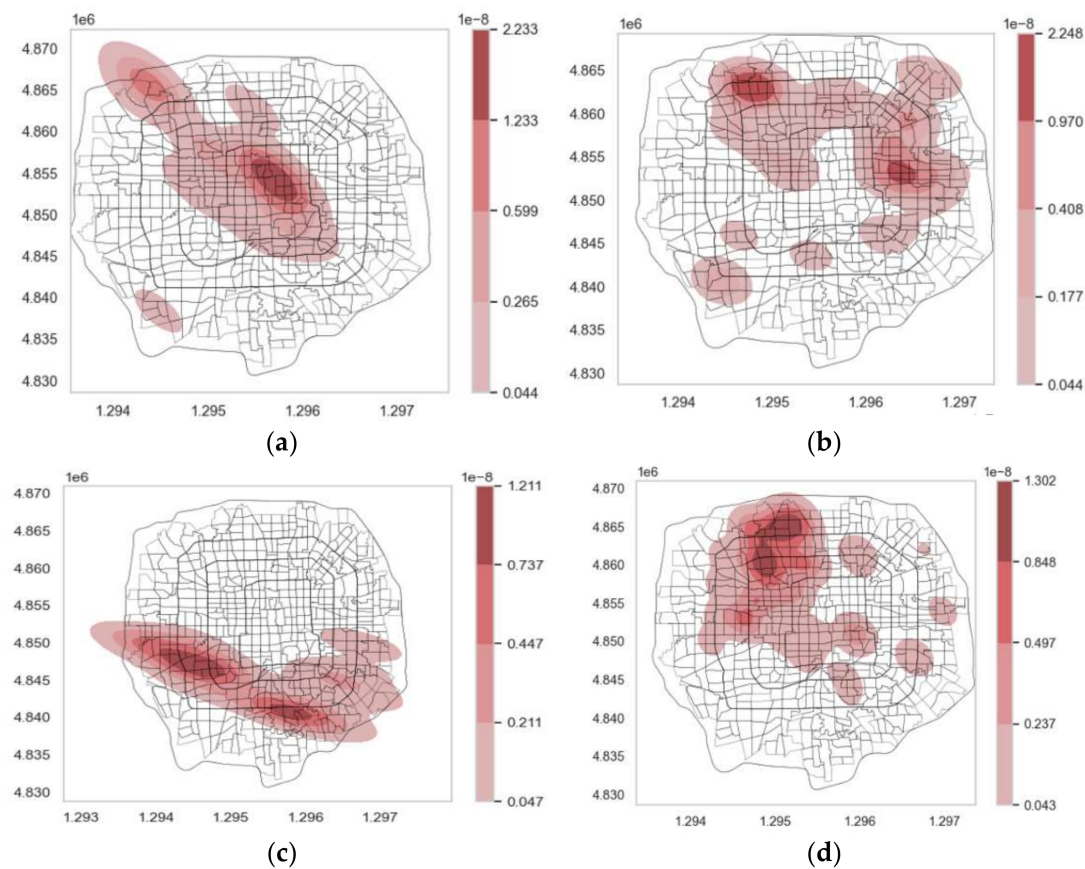


Figure 7. Intensity of representative kinds of POIs. (a) Intensity of toponymic address information, historical sites, scenic spots, cultural relics, and aquariums in cluster C1; (b) Intensity of business trade, corporate enterprises, restaurants, and bars in cluster C2; (c) Intensity of factories, industries, clothing, shoes, hats, and leather goods stores in cluster C7; (d) Intensity of science, education and cultural services, schools, and educational training institutions in cluster C8.

Developed Residential Areas (C3)

C3 was annotated as a developed residential area based on the following two aspects. First, the strong interaction between C3 and C2 indicates that C3 possesses the properties of a residential area (Figure 6a–d). Most residents leave from C3 to C2 during the morning (8:00–9:00) and leave from C2 to C3 during the evening peak (17:00–19:00) on weekdays; however, such a regular commuting pattern does not exist on the weekends. Second, Table 2 indicates that C3 has the maximum number of dwellings, accommodation services, courier services, living services, hair dressings, healthcare treatment, banks, and convenience stores, which constitute a mature and complete POI configuration of developed residential areas.

Emerging Residential Areas (C4)

Figure 6e–h indicates that a workday commuting pattern similar to C3 also exists in C4, implying that C4 also has residential area characteristics. Table 2 indicates that C4 has a POI configuration similar to that of C3, but the frequency density and category rate of the POIs associated with living services in C4 are relatively smaller; it is worth mentioning that C4 has more moving companies and intermediary institutions than C3, which are representative POIs of emerging residential areas; moreover, the regions in C4 are primarily distributed in the Chaoyang District, with several emerging business circles (Wangjing

Street and Yansha Center) and newly developing dwellings. Therefore, C4 was annotated as an emerging residential area.

Less Developed Residential Areas (C5)

C5 was annotated as a less-developed residential area. First, Figure 6i–l shows that C5 also presented a commuting pattern with C2 on weekdays, indicating that C5 had the attributes of residential areas. Second, Table 2 shows that the dominant POIs related to living services in C5 are sufficient to provide normal living needs for residents in all aspects, but they are lower than those in C3 and C4. The typical POI located in C5 is house number information, commonly known as “doorplate”, which is the symbol of some ancient architecture in Beijing, such as “hutong” and “quadrangle courtyards”.

Residential/Commercial/Entertainment Areas (C6)

Table 2 shows that the POI configuration in C6 is the most complicated among the clusters. The diversity of the POIs determines C6 as a mix of residential, commercial, and entertainment areas. Living service, healthcare treatments, convenience stores, express service, shopping malls, buildings, companies, enterprises, theaters, attractions, and recreation and amusement places are relatively sufficient and balanced in this cluster, which illustrates that C6 has the characteristics of three kinds of land-use types. The workday commuting pattern shown in Figure 6m–p also indicates that C6 has the properties of a residential area. In addition, the frequency density and category rate of companies and enterprises, shopping malls, buildings, hotels, theaters, and catering services are second only to those in C2, which illustrates that C6 has strong commercial attributes; moreover, a number of well-known historic sites and entertainment places are located in C6, such as Maodun’s former residence, Memorial Hall, South Luogu Alley, Prince Gong’s Mansion, and the Yuhe ancient road site. Therefore, C6 was annotated as a mix of residential, commercial, and entertainment areas.

Industrial Areas (C7)

C7 was annotated as an industrial area as the representative POIs in this cluster were clothing factories and industries. The frequency density and category rate of these two POIs in this cluster were the highest among the ten clusters (Table 2). Figure 7c visualizes the intensity of several representative types of POIs, such as factories, industries, clothing, shoes, hats, and leather goods stores. Regions in C7 are mainly distributed in the Fengtai District, which contain clothing trade markets, shoes, hats, and leather goods manufacturers, and large technology parks.

Education and Scientific Research Areas (C8)

C8 was annotated as an education and scientific research area as the frequency density and category rate of educational culture services and scientific institution schools in C8 were the highest among the ten clusters (Table 2). Figure 7d visualizes the intensity of several representative types of POIs, such as science, education and cultural services, schools, and educational training institutions; it was found that major scientific and educational sites in Beijing are located in this cluster, such as Beijing University, Tsinghua University, Chinese Academy of Sciences, Beijing University of Aeronautics and Astronautics, Beijing University of science and technology, Beijing Jiaotong University, and Beijing Normal University.

Green Space Areas (C9)

C9 was annotated as a greenland area; it can be found from Table 2 that the type and number of POIs in C9 are the least compared to those in the other nine clusters. The only representative POI in this cluster is the name of the natural place, which is consistent with public green spaces, mountains, and parks such as Shijing Mountain, Babao Mountain, Beijing Olympic Park, Bihai Park, and Guta Park.

Road and Transportation Areas (C10)

C10 was annotated as a road and transportation area as the frequency density and category rate of transportation services, road ancillary facilities, Sionpec, gas stations, long-distance bus stations, and railway stations were the highest among the ten clusters (Table 2). Beijing's main large-scale train stations and passenger transportation stations are concentrated in this cluster, such as the Beijing Railway Station, Liuliqiao passenger transport hub, and Yongdingmen long-distance station.

4.2.4. Quantitative Comparison and Analysis

In Figure 8, the clustering results of the four methods are displayed. The land-use types of the clusters obtained by the comparative methods were annotated using the same strategies described in Section 4.2.3. Similar to the method used by Pei et al. [3], the land-use types of traffic analysis zones identified by different methods were compared with the Beijing land-use planning map (Figure 4b). Table 3 lists the overall accuracy of the land-use classification results identified using different methods. Overall accuracy was calculated as follows:

$$OA = \frac{\sum_{j=1}^{N_{TAZ}} I_j}{N_{TAZ}} \times 100\%, \quad (12)$$

where N_{TAZ} is the number of traffic analysis zones. I is an indicator function, if the land use type of a traffic analysis zone is consistent with that in the land-use planning map, then $I = 1$; else, $I = 0$.

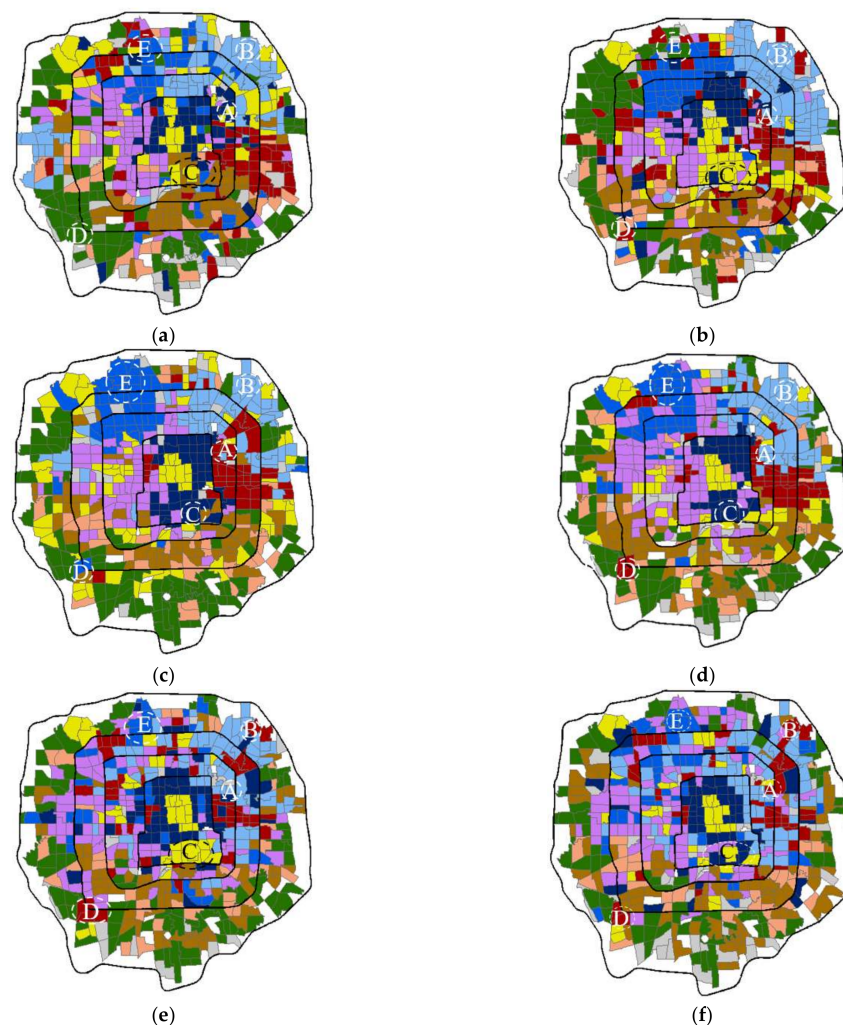


Figure 8. Cont.

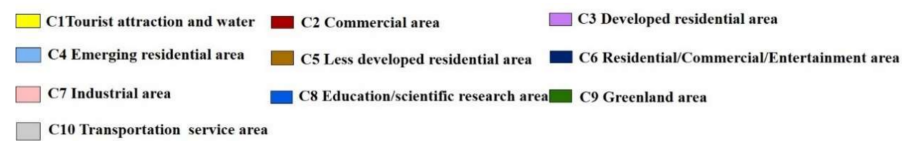


Figure 8. Clustering results by using different methods. (a) Single-view spectral clustering method using taxi GPS trajectory; (b) Single-view spectral clustering method using bus smart card transaction data; (c) Single-view spectral clustering method using POI data; (d) Weighted average spectral clustering method; (e) gLMSC; (f) agc2msc.

Table 3. Overall accuracy of different methods.

| Methods | Single-View Method (Taxi) | Single-View Method (Bus) | Single-View Method (POI) | Weighted Average Method | gLMSC | agc2msc |
|------------------|---------------------------|--------------------------|--------------------------|-------------------------|--------|---------|
| Overall accuracy | 44.53% | 45.93% | 50.95% | 53.15% | 64.82% | 68.11% |

agc2msc achieved the highest classification accuracy of 68.11%. Examples are provided to illustrate the advantages of the proposed method. Some commercial areas were not identified using the single-view clustering method, such as Sanlitun in region A, Wangjing Street in region B, and Advanced Business Park in region D (Figure 8a–c). In Figure 8a,b, although Temple of Heaven in region C could be accurately identified using only taxi GPS trajectories or bus smart-card transaction data, its surrounding areas were all misclassified. In Figure 8c, Temple of Heaven in region C and China University of Mining and Technology, Chinese Academy of Sciences, and China Agricultural University in region E were overestimated when using POI data. The weighted average clustering method could not identify some commercial areas in regions A and B and the Temple of Heaven in region C (Figure 8d). Additionally, the scientific research and education areas in region E were overestimated. In Figure 4e, Sanlitun in region A was misclassified, and Temple of Heaven in region C and Advanced Business Park in region D were overestimated. From Figure 8f, we found that all the misclassified areas in Figure 8a–e could be correctly identified using the method proposed in this study.

5. Discussion

The experimental results on benchmark and real-world datasets show that agc2msc performs better than the existing single-view, weighted average, and multi-view subspace clustering methods. The primary reasons for this are as follows:

- Compared with the single-view spectral clustering method, the complementary information of multi-source geospatial big data can be incorporated accurately using agc2msc. Therefore, agc2msc can alleviate the bias problem caused by a single type of geospatial big data and comprehensively describe urban structures and organizations in cities.
- Compared with the weighted average spectral clustering strategy, agc2msc can fuse the shared and complementary information among different types of geospatial big data. The underlying structure of the multi-source data can be reconstructed accurately using agc2msc. Therefore, agc2msc can capture the complementarity of multi-source geospatial big data more accurately.
- Compared with the multi-view subspace clustering method, agc2msc can construct appropriate neighboring relationships for high-dimensional, noisy, and non-uniform geospatial big data. A more robust shared subspace can be obtained under the constraint of a shared nearest neighbor graph.

The clusters detected from multi-source geospatial big data using agc2msc may provide potential application value for urban planning.

- (i) Actually, these clusters can reveal urban function zones in a city. By fusing multi-source geospatial big data, we can obtain a comprehensive view of urban function zones naturally formulated according to human activities. The clusters identified by agc2msc may be further used for public services, business site selection, and human-centric urban planning [30]. The more complex land-use types identified in this study can also provide a reference for urban planning and city development. For instance, residential areas can be divided into three types: developed residential areas; emerging residential areas; less developed residential areas; and a mix of residential, commercial, and entertainment areas. Scientific and research areas were also identified. Existing remote sensing techniques are hard to obtain this complex division. The clustering results obtained by agc2msc may help urban planners make more strategic decisions and improve the quality of land-use mapping.
- (ii) Some calibrations may be presented for urban land-use planning. Although the detection rate of the proposed method is relatively low (overall accuracy is 68.11%), the clustering results are useful for infer the actual land use which cannot be captured by land use planning map. In fact, the actual land use types may differ from the Beijing land use planning map (Figure 4b). The land use planning map was obtained based on the physical characteristics of ground components (e.g., spectral, shape, and texture); therefore, the land use planning map is hard to reflect the actual way of how people use spaces [3,26]. We give some examples to illustrate that the actual land use types of some traffic analysis zones are not consistent with those in Beijing land-use planning map.

In Figure 4b, area A and its surrounding areas are labeled as tourist attractions and water land, area B is labeled as a mixture of tourist attractions and residential areas, and areas C and D are labeled as industrial land and green land, respectively; however, our method could distinguish area A from its surrounding areas, areas B and C were annotated as commercial land, and area D was annotated as tourist attraction land. In Figure 9, although area A is a well-known scenic spot (Temple of Heaven), its surrounding areas do not contain any tourist attractions, but are a mixture of residential and entertainment commercial areas. Area B is Sanlitun, one of the most prosperous commercial areas and amusement streets for nightlife in Beijing. The landmark of area C is Advanced Business Park, which is a Sino-foreign joint venture project with the largest single area since the opening of Zhongguancun. In recent years, Advanced Business Park has developed into an emerging economic experiment zone and prosperous business district in Beijing, including a large number of hotels, financial squares, and underground commercial streets. Area D is covered by Beijing World Park and Huake Golf Club. Beijing World Park has been rated as a Beijing 4A scenic spot. Therefore, area D was annotated as a tourist attraction area.

Based on the above analysis, we argue that the main reason of the low detection rate of agc2msc may be the mismatch between the physical characteristics and social function of urban land. The urban land use types identified by the proposed method can provide valuable calibration and reference for urban planning. In addition, the accurate of the clustering results will also be influenced by some limitations of the proposed method (details can be found in Section 6).

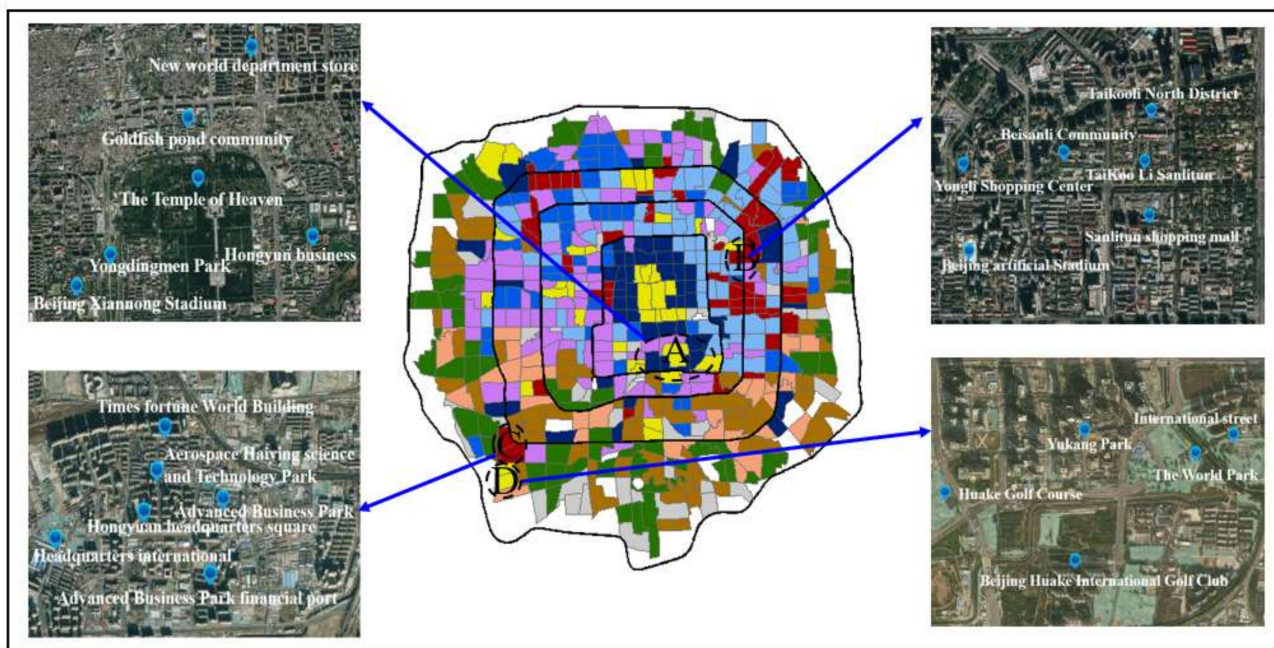


Figure 9. Urban land use types of some regions identified using Google Earth images.

6. Conclusions

This study developed a method with adaptive graphs to constrain multi-view subspace clustering of geospatial big data from multiple sources. We regarded multi-source geospatial big data to describe urban structures and used a multi-view learning strategy to fuse the information embedded in different sources of geospatial data. Therefore, the bias of a single type of data could be reduced. The neighboring relationships of high-dimensional, noisy, and non-uniform geospatial big data were appropriately constructed using the shared nearest neighbor graph. The graph was used as a constraint to obtain a more robust subspace shared by multi-source geospatial big data. Experiments on benchmark datasets and multi-source geospatial big data in Beijing showed that agc2msc outperforms the typical single-view, weighted average, and multi-view subspace clustering methods. The urban land use inferred by the proposed method may provide a useful calibration and reference for urban land-use planning.

Although agc2msc provides a powerful tool for clustering multi-source geospatial big data, it has three limitations. First, we only extracted temporal rhythm features from human mobility data, and some complex information (e.g., the interactions between different traffic analysis zones) may be neglected. Second, the proposed method does not distinguish shared information among multiple views and the view-specific information of each view. In the future, the latent representation of each type of data should be segregated into shared and specific parts. Third, although traffic analysis zones are reliable spatial units in urban studies, the effect of the modifiable areal unit problem on geospatial big data clustering should be further analyzed.

Author Contributions: Conceptualization, Q.L. and M.D.; methodology, Q.L. and W.H.; software, W.H.; validation, W.H. and Q.L.; formal analysis, W.H. and Q.L.; investigation, Q.L., Weihua Huan, and M.D.; resources, Q.L. and W.H.; data curation, W.H.; writing—original draft preparation, W.H.; writing—review and editing, Q.L. and M.D.; visualization, W.H.; supervision, M.D.; project administration, Q.L. and M.D.; funding acquisition, Q.L. and M.D. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded through support from the National Natural Science Foundation of China (NSFC) (No. 41971353 and 41730105), the Natural Science Foundation of Hunan Province (No. 2021JJ20058 and 2020JJ4695), and Water conservancy science and technology project of Guizhou Province (No. KT202110 and KT202002).

Data Availability Statement: Not applicable.

Acknowledgments: We would like to thank Tao Pei and Zhou Huang for providing the land use data and traffic analysis zone data in Beijing. The authors gratefully acknowledge the comments from the reviewers.

Conflicts of Interest: The authors declare no conflict of interest.

Appendix A. Solution and Optimization of the Objective Function

The model was pretrained based on the Adam optimizer [58]. The objective function can be transformed into the form as follows [23]:

$$L = \min_{f^{(M,k)}, Z} \frac{1}{2} \sum_{v=1}^V \sum_{i=1}^N \left(\|x_i^{(v)} - f_i^{(M,v)}\|^2 + \alpha \|f_i^{(\frac{M}{2},v)} - f_i^{(\frac{M}{2},v)} Z\|^2 \right) + \beta \text{Tr}(Z^T L Z), \quad (\text{A1})$$

where $\hat{X}^{(v)} = [f_1^{(M,v)}, f_2^{(M,v)}, \dots, f_N^{(M,v)}]$, $H^{(v)} = [f_1^{(M/2,v)}, f_2^{(M/2,v)}, \dots, f_N^{(M/2,v)}]$, $f_i^{(m,v)}$ is the output of the m^{th} layer of the autoencoder networks of sample x_i in the v^{th} view. As the entry of Equation (A1), $f_i^{(m,v)}$ is calculated by the activation function:

$$f_i^{(m,v)} = g(W^{(m,v)} f_{i-1}^{(m-1,v)} + b^{(m,v)}), \quad (\text{A2})$$

where $W^{(m,v)}$ and $b^{(m,v)}$ denote the weight and bias of the m^{th} layer of the auto-coder networks in the v^{th} view, respectively.

The gradient of the loss function L in Equation (A1) with $W^{(m,v)}$, $b^{(m,v)}$, Z are given as follows [22,62]:

$$\frac{\partial L}{\partial W^{(m,v)}} = \sum_{i=1}^N \left(\Delta_i^{(m,v)} + \alpha \Lambda_i^{(m,v)} \right) \left(f_i^{(m-1,v)} \right)^T, \quad (\text{A3})$$

$$\frac{\partial L}{\partial b^{(m,v)}} = \sum_{i=1}^N \Delta_i^{(m,v)} + \alpha \Lambda_i^{(m,v)}, \quad (\text{A4})$$

$$\frac{\partial L}{\partial Z} = \alpha \left(H^{(k)T} H^{(k)} Z - H^{(k)T} H^{(k)} \right) + \beta \left(L^T + L \right) Z, \quad (\text{A5})$$

where $\Delta_i^{(m,v)}$ and $\Lambda_i^{(m,v)}$ are defined as:

$$\Delta_i^{(m,v)} = \begin{cases} -\left(x_i^{(v)} - f_i^{(m,v)}\right) \odot g'\left(y_i^{(m,v)}\right), m = M \\ \left(W^{(m+1,v)}\right)^T \Delta_i^{(m+1,v)} \odot g'\left(y_i^{(m,v)}\right), m < M \end{cases}, \quad (\text{A6})$$

$$\Lambda_i^{(m,v)} = \begin{cases} 0, m \geq \frac{M}{2} + 1 \\ \left(f_i^{(m,v)} - z_i - f_i^{(m,v)} z_i^T z_i\right) \odot g'\left(y_i^{(m,v)}\right), m = \frac{M}{2}, \\ \left(W^{(m+1,v)}\right)^T \Lambda_i^{(m+1,v)} \odot g'\left(y_i^{(m,v)}\right), m \leq \frac{M}{2} - 1 \end{cases}, \quad (\text{A7})$$

where z_i denotes the i^{th} column of Z , $y_i^{(m,v)} = W^{(m,v)} f_i^{(m-1,v)} + b^{(m,v)}$.

Appendix B. Experiments on Benchmark Datasets

Appendix B.1. Benchmark Datasets

The detailed statistical information of the four datasets was given in Table A1. In each dataset, the cluster label of each sample is available [22,39].

Table A1. Statistical information of the four multi-view datasets.

| Datasets | Samples | Clusters | Views | Dimensions |
|------------|---------|----------|-------|-------------------------|
| ORL | 400 | 40 | 3 | 4096/3304/6750 |
| Yale | 165 | 15 | 3 | 4096/3304/6750 |
| MSRCV1 | 210 | 7 | 6 | 1032/48/512/100/256/210 |
| Caltech101 | 441 | 7 | 6 | 1032/32/512/64/256/441 |

- (i) ORL: This dataset contains 400 gray images of 40 different individuals; it was created by Olivetti research laboratory in Cambridge, UK from April 1992 to April 1994. The dataset contains a total of 40 directories. Each directory represents 10 facial pictures taken by the same person at different times and in different environments (e.g., light/no light, glasses/no glasses, changes in different facial expressions, etc.). All the pictures are stored in the form of gray-scale image and Portable Gray Map format, and the picture size is 92×112 . Three types of features were used in the experiment, i.e., intensity (4096 dimensions), Local Binary Pattern (3304 dimensions) and Gabor (6750 dimensions), representing three different views of observation.
- (ii) Yale: This dataset was created by Yale University and contains a total of 165 grayscale images from 15 different individuals. The dataset contains a total of 15 directories. Each directory represents 11 face images of the same person under different expressions, gestures and illumination. The size of each image is 100×100 . Variations of images include central light/edge light, wearing glasses/not wearing glasses, happiness/sadness, surprise/blink, etc. Similar to ORL dataset, three different types of features were extracted in the experiment as three different views of observation, namely intensity (4096 dimension), Local Binary Pattern (3304 dimension) and Gabor (6750 dimension).
- (iii) MSRCV1: This dataset contains 210 images from 7 categories collected from 6 different views, and each category contains 30 images. From the collected samples, it can be seen that seven categories include human face, animals, trees, scenery, bicycles, cars, planes. Six types of high-dimensional features were extracted in the experiment: Centrist (view1), Charcot Marie Tooth (view2), Gist (view3), Histogram of Oriented Gradient (view4), Local Binary Pattern (view5), and Scale-invariant feature transform (view6).
- (iv) Caltech101-7: Caltech101 is a dataset widely used in image classification in deep learning, which contains 101 types of images. In this study, the subset of this dataset (i.e., Caltech101-7) was used in the experiment. In Caltech101-7, a total of 441 images of 7 categories were selected, including face, coin, Garfield, motorcycle, Snoopy, parking sign and chair. Six types of high-dimensional features were extracted for experiment, which was similar to those of MSRCV1.

Appendix B.2. Baselines

To demonstrate the effectivity of agc2msc, two single-view clustering methods and seven state-of-the-art multi-view clustering methods were compared.

- (1) $Single_{best}$: This method performs the standard spectral clustering [28] on the most informative view.
- (2) LRR_{best} : This method performs single-view algorithm LRR [34] on the most informative view.
- (3) FeatConcat [46]: This method directly concatenates the features from different views, and then applies the concatenated features to single-view clustering algorithm.
- (4) RMSC [63]: This method firstly recovers a shared low-rank transition probability matrix, and then uses a Markov chain to cluster.
- (5) gLMS [18]: This method firstly calculates an underlying latent representation shared by multi-view features, and then applies the latent representation to subspace clustering.

- (6) DiMSC [38]: This method extends the existing single-view subspace clustering into multi-view domain and exploits the complementary information of multi-view representations by enforcing Hilbert Schmidt Independence Criterion term.
- (7) CSMSC [39]: This method exploits both the consistent and specific information among multi-view features by pursuing a view-consistent representation matrix and a set of view-specific self-representation matrices.
- (8) DSS-MSC [22]: This method decomposes the underlying latent representation into shared component and view-specific components, which exploit the underlying correlations cross multiple views and simultaneously capture specific property for each independent view.
- (9) MSCNLG [23]: This method introduces artificial neural network under each view to obtain a set of latent representations and integrates local and global graph information into self-expressive layers.

Appendix B.3. Evaluation Metrics and Clustering Results

To assess the performance of agc2msc and the baseline methods, six clustering evaluation indexes (NMI, ACC, F-score, AR, precision and recall) were used to measure the quality of clustering results. In Tables A2–A5, clustering results on four multi-view benchmark datasets were reported; it can be found that agc2msc in this study outperforms the comparative methods.

Table A2. Clustering results of the ten methods on ORL.

| Method | NMI | ACC | F-Score | AR | Precision | Recall |
|------------------------------|---------------|---------------|---------------|---------------|---------------|---------------|
| <i>Single_{best}</i> | 0.884 ± 0.002 | 0.726 ± 0.025 | 0.664 ± 0.005 | 0.655 ± 0.005 | 0.610 ± 0.006 | 0.728 ± 0.005 |
| <i>LRR_{best}</i> | 0.895 ± 0.006 | 0.773 ± 0.003 | 0.731 ± 0.004 | 0.724 ± 0.020 | 0.701 ± 0.001 | 0.754 ± 0.002 |
| FeatConcat | 0.831 ± 0.003 | 0.648 ± 0.033 | 0.564 ± 0.007 | 0.553 ± 0.007 | 0.522 ± 0.007 | 0.614 ± 0.008 |
| RMSC | 0.872 ± 0.012 | 0.723 ± 0.025 | 0.654 ± 0.028 | 0.645 ± 0.029 | 0.607 ± 0.033 | 0.709 ± 0.027 |
| gLMSC | 0.924 ± 0.011 | 0.830 ± 0.017 | 0.771 ± 0.028 | 0.765 ± 0.044 | 0.728 ± 0.010 | 0.819 ± 0.010 |
| DiMSC | 0.940 ± 0.003 | 0.838 ± 0.001 | 0.807 ± 0.003 | 0.802 ± 0.000 | 0.764 ± 0.012 | 0.856 ± 0.004 |
| CSMSC | 0.942 ± 0.005 | 0.868 ± 0.012 | 0.831 ± 0.001 | 0.615 ± 0.005 | 0.673 ± 0.002 | 0.610 ± 0.006 |
| DSS-MSC | 0.928 ± 0.010 | 0.795 ± 0.010 | 0.766 ± 0.010 | 0.762 ± 0.010 | 0.719 ± 0.010 | 0.823 ± 0.010 |
| MSCNLG | 0.936 ± 0.002 | 0.885 ± 0.003 | 0.857 ± 0.004 | 0.825 ± 0.002 | 0.885 ± 0.002 | 0.885 ± 0.002 |
| agc2msc | 0.943 ± 0.002 | 0.893 ± 0.002 | 0.871 ± 0.002 | 0.831 ± 0.002 | 0.890 ± 0.002 | 0.890 ± 0.002 |

Table A3. Clustering results of the ten methods on Yale.

| Method | NMI | ACC | F-Score | AR | Precision | Recall |
|------------------------------|---------------|---------------|---------------|---------------|---------------|---------------|
| <i>Single_{best}</i> | 0.654 ± 0.009 | 0.616 ± 0.030 | 0.475 ± 0.043 | 0.440 ± 0.011 | 0.457 ± 0.011 | 0.495 ± 0.010 |
| <i>LRR_{best}</i> | 0.709 ± 0.011 | 0.697 ± 0.001 | 0.547 ± 0.007 | 0.515 ± 0.004 | 0.529 ± 0.003 | 0.567 ± 0.004 |
| FeatConcat | 0.648 ± 0.030 | 0.607 ± 0.043 | 0.471 ± 0.039 | 0.434 ± 0.042 | 0.447 ± 0.043 | 0.497 ± 0.032 |
| RMSC | 0.872 ± 0.012 | 0.723 ± 0.025 | 0.654 ± 0.028 | / | / | / |
| gLMSC | 0.735 ± 0.021 | 0.752 ± 0.026 | 0.564 ± 0.019 | 0.551 ± 0.011 | 0.543 ± 0.015 | 0.571 ± 0.013 |
| DiMSC | 0.727 ± 0.010 | 0.709 ± 0.003 | 0.564 ± 0.002 | 0.535 ± 0.001 | 0.543 ± 0.001 | 0.586 ± 0.003 |
| CSMSC | 0.784 ± 0.001 | 0.752 ± 0.007 | 0.640 ± 0.004 | 0.615 ± 0.005 | 0.673 ± 0.002 | 0.610 ± 0.006 |
| DSS-MSC | 0.779 ± 0.021 | 0.782 ± 0.013 | 0.613 ± 0.012 | 0.601 ± 0.009 | 0.529 ± 0.010 | 0.622 ± 0.015 |
| MSCNLG | 0.879 ± 0.002 | 0.903 ± 0.002 | 0.831 ± 0.002 | 0.790 ± 0.002 | 0.903 ± 0.002 | 0.903 ± 0.002 |
| agc2msc | 0.913 ± 0.002 | 0.925 ± 0.002 | 0.871 ± 0.002 | 0.835 ± 0.002 | 0.922 ± 0.002 | 0.922 ± 0.002 |

Table A4. Clustering results of the ten methods on MSRCV1.

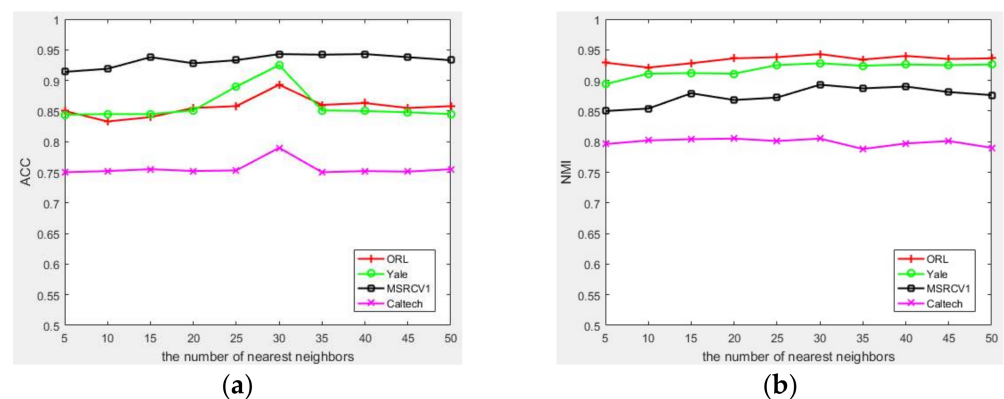
| Method | NMI | ACC | F-Score | AR | Precision | Recall |
|-----------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| $Single_{best}$ | 0.574 ± 0.032 | 0.668 ± 0.051 | 0.535 ± 0.043 | 0.536 ± 0.010 | 0.571 ± 0.009 | 0.612 ± 0.009 |
| LRR_{best} | 0.569 ± 0.008 | 0.676 ± 0.009 | 0.524 ± 0.009 | 0.502 ± 0.010 | 0.543 ± 0.009 | 0.587 ± 0.007 |
| FeatConcate | 0.613 ± 0.042 | 0.672 ± 0.031 | 0.575 ± 0.024 | 0.505 ± 0.032 | 0.566 ± 0.021 | 0.586 ± 0.027 |
| RMSC | 0.650 ± 0.022 | 0.750 ± 0.048 | 0.628 ± 0.023 | / | / | / |
| gLMSC | 0.752 ± 0.011 | 0.848 ± 0.013 | 0.738 ± 0.018 | 0.721 ± 0.017 | 0.744 ± 0.012 | 0.743 ± 0.011 |
| DiMSC | 0.692 ± 0.002 | 0.810 ± 0.002 | 0.685 ± 0.002 | 0.634 ± 0.002 | 0.679 ± 0.002 | 0.691 ± 0.002 |
| CSMSC | 0.756 ± 0.002 | 0.857 ± 0.002 | 0.756 ± 0.002 | 0.717 ± 0.002 | 0.750 ± 0.002 | 0.762 ± 0.002 |
| DSS-MSC | 0.743 ± 0.015 | 0.846 ± 0.011 | 0.726 ± 0.021 | 0.681 ± 0.014 | 0.711 ± 0.011 | 0.743 ± 0.013 |
| MSCNLG | 0.850 ± 0.002 | 0.921 ± 0.002 | 0.862 ± 0.002 | 0.830 ± 0.002 | 0.922 ± 0.002 | 0.922 ± 0.002 |
| agc2msc | 0.893 ± 0.002 | 0.943 ± 0.002 | 0.900 ± 0.002 | 0.869 ± 0.002 | 0.945 ± 0.002 | 0.942 ± 0.002 |

Table A5. Clustering results of the ten methods on Caltech101.

| Method | NMI | ACC | F-Score | AR | Precision | Recall |
|-----------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| $Single_{best}$ | 0.589 ± 0.009 | 0.629 ± 0.007 | 0.576 ± 0.009 | 0.523 ± 0.012 | 0.586 ± 0.014 | 0.566 ± 0.003 |
| LRR_{best} | 0.639 ± 0.002 | 0.646 ± 0.003 | 0.649 ± 0.002 | 0.580 ± 0.001 | 0.631 ± 0.001 | 0.623 ± 0.003 |
| FeatConcate | 0.603 ± 0.017 | 0.641 ± 0.020 | 0.601 ± 0.023 | 0.526 ± 0.034 | 0.624 ± 0.021 | 0.579 ± 0.024 |
| gLMSC | 0.694 ± 0.013 | 0.722 ± 0.012 | 0.683 ± 0.009 | 0.620 ± 0.002 | 0.670 ± 0.002 | 0.695 ± 0.002 |
| DiMSC | 0.679 ± 0.002 | 0.746 ± 0.002 | 0.709 ± 0.002 | 0.653 ± 0.002 | 0.717 ± 0.002 | 0.702 ± 0.002 |
| CSMSC | 0.701 ± 0.002 | 0.732 ± 0.002 | 0.702 ± 0.002 | 0.630 ± 0.002 | 0.680 ± 0.002 | 0.702 ± 0.002 |
| DSS-MSC | 0.691 ± 0.002 | 0.737 ± 0.001 | 0.703 ± 0.006 | 0.635 ± 0.002 | 0.698 ± 0.002 | 0.710 ± 0.002 |
| MSCNLG | 0.758 ± 0.002 | 0.764 ± 0.002 | 0.760 ± 0.002 | 0.687 ± 0.002 | 0.748 ± 0.002 | 0.748 ± 0.002 |
| agc2msc | 0.805 ± 0.002 | 0.790 ± 0.002 | 0.793 ± 0.002 | 0.710 ± 0.002 | 0.762 ± 0.002 | 0.762 ± 0.002 |

Appendix B.4. Parameter Settings and Parameter Sensitivity

We evaluated the variation of ACC and NMI on four benchmark datasets with different numbers of nearest neighbors. In Figure A1, when the number of nearest neighbors was set to 30, both ACC and NMI reached the maximum value. Therefore, the number of nearest neighbors was set to 30 in the experiment.

**Figure A1.** Variation of ACC and NMI with different numbers of nearest neighbors on four benchmark datasets, (a) ACC, (b) NMI.

In addition, parameter sensitivity on balancing the self-representation term (alpha) and the adaptive graph constraint term (beta), can be shown as Figures A2–A5. Taking Figure A2a as an example, it can be found that when $\alpha \in [0.001, 0.01]$, $\beta \in [0.001, 10]$, the values of NMI, ACC, F-score, and AR are all stable and high. Therefore, we can determine the range of alpha and beta. The interpretation of Figures A3–A5 is similar to that of Figure A2.

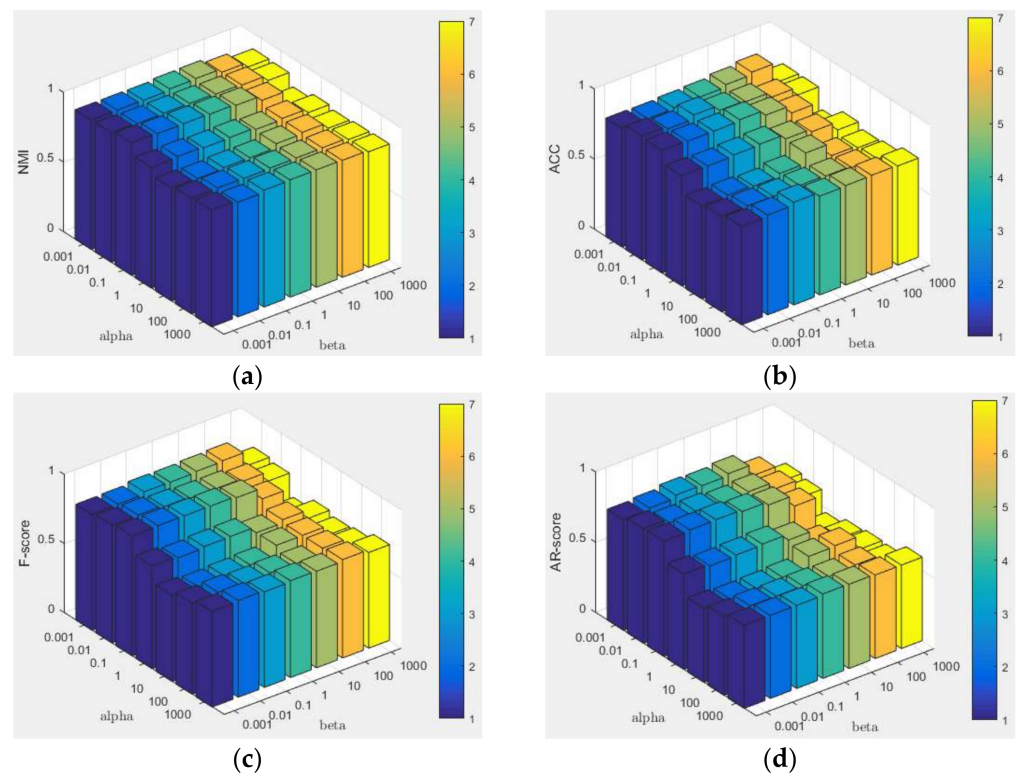


Figure A2. Sensitivity test on alpha and beta versus four metrics on ORL. (a) NMI, (b) ACC, (c) F-score, (d) AR.

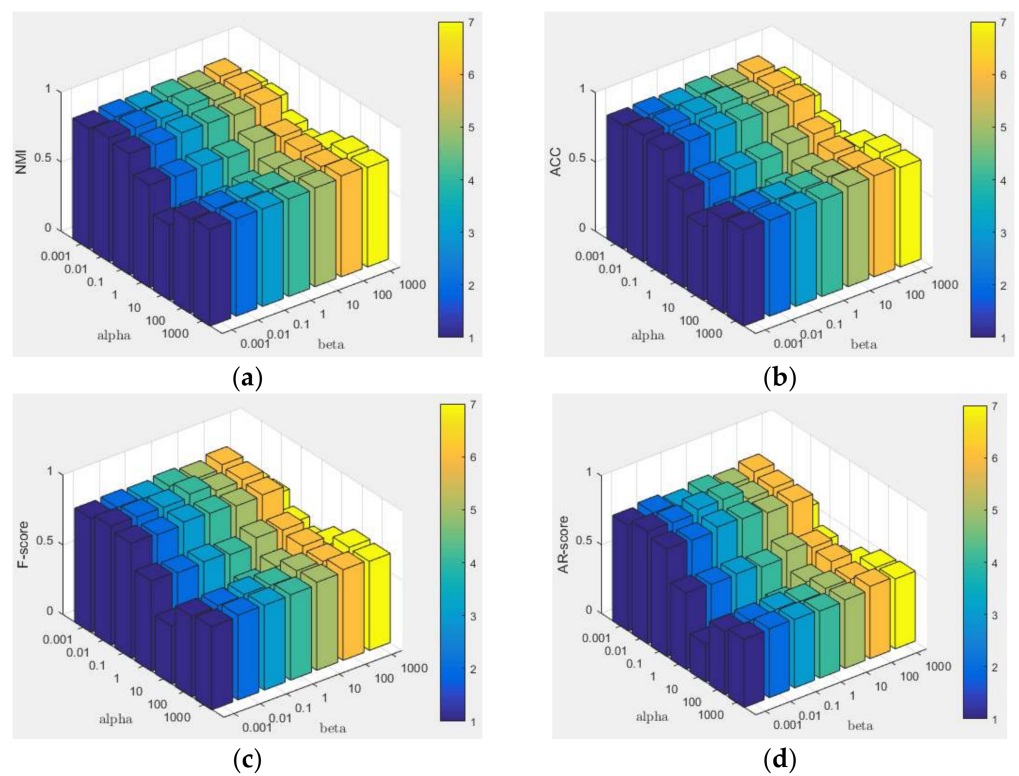


Figure A3. Sensitivity test on alpha and beta versus four metrics on Yale. (a) NMI, (b) ACC, (c) F-score, (d) AR.

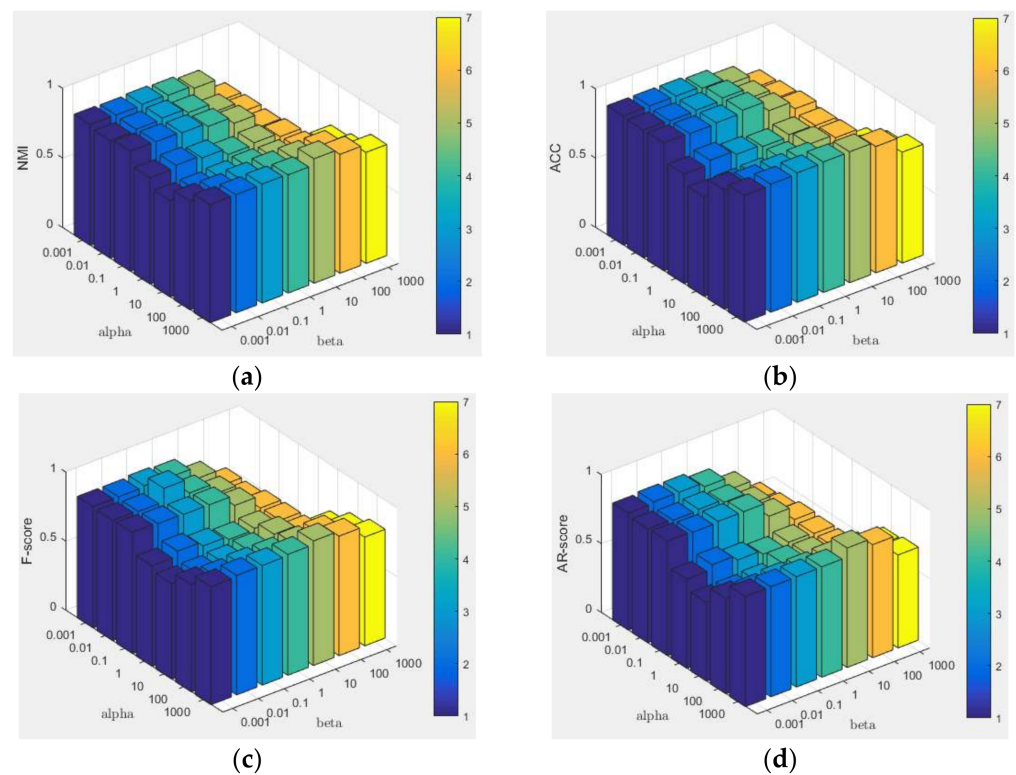


Figure A4. Sensitivity test on alpha and beta versus four metrics on MSRCV1. (a) NMI, (b) ACC, (c) F-score, (d) AR.

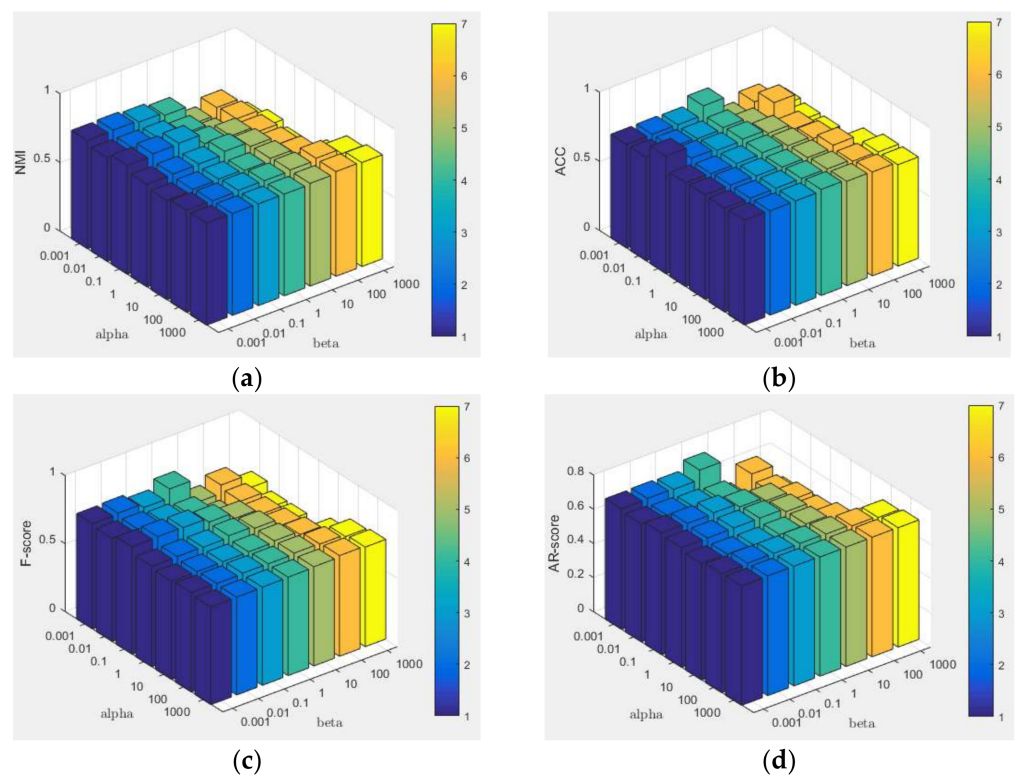


Figure A5. Sensitivity test on alpha and beta versus four metrics on Caltech101. (a) NMI, (b) ACC, (c) F-score, (d) AR.

Appendix B.5. Convergence Analysis

The convergence of the proposed method on four benchmark datasets was presented in Figure A6. When epochs reach 0.2×10^3 , the proposed method shows a stable convergence on each dataset.

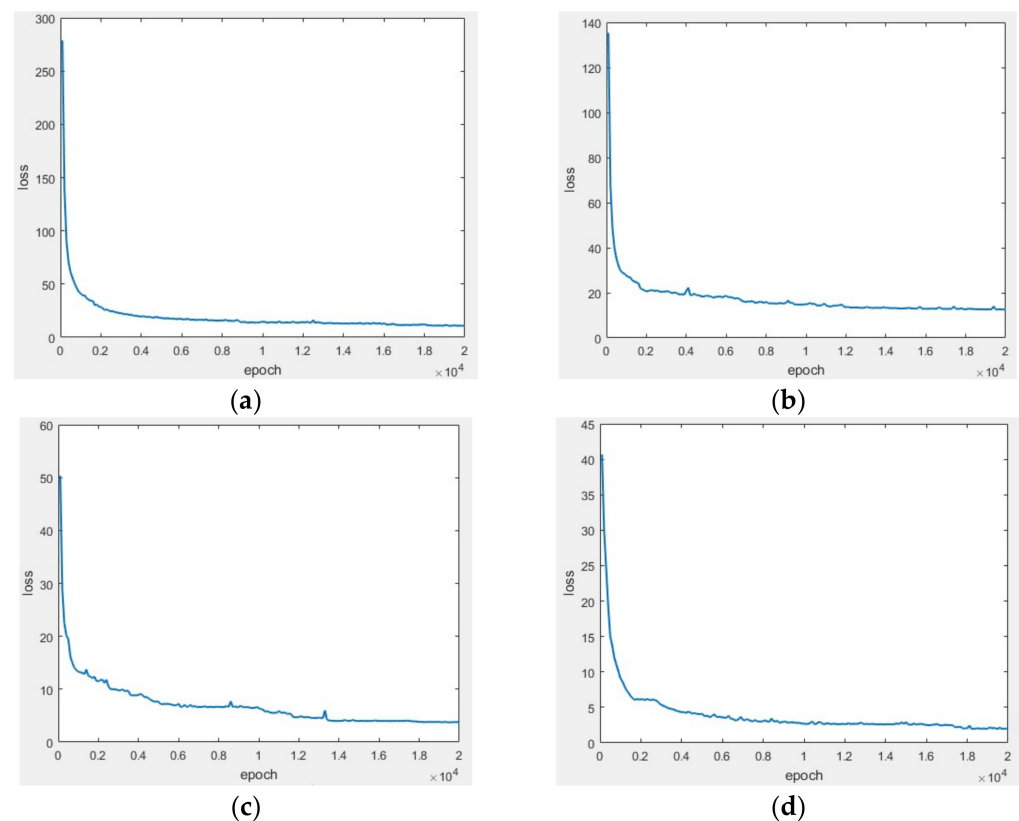


Figure A6. Convergence curves with training epochs. (a)ORL, (b)Yale, (c)MSRCV1, (d)Caltech101.

References

- Pan, G.; Qi, G.; Wu, Z.; Zhang, D.; Li, S. Land-use classification using taxi gps traces. *IEEE Trans. Intell. Transp. Syst.* **2013**, *14*, 113–123. [\[CrossRef\]](#)
- Long, Y.; Shen, Z. Discovering functional zones using bus smart card data and points of interest in Beijing. In *Geospatial Analysis to Support Urban Planning in Beijing*; Long, Y., Shen, Z., Eds.; Springer: Berlin, Germany, 2015; Volume 116, pp. 193–217.
- Pei, T.; Sobolevsky, S.; Ratti, C.; Shaw, S.L.; Li, T.; Zhou, C. A new insight into land use classification based on aggregated mobile phone data. *Int. J. Geogr. Inf. Sci.* **2014**, *28*, 1988–2007. [\[CrossRef\]](#)
- Comito, C.; Pizzuti, C.; Procopio, N. Online clustering for topic detection in social data streams. In Proceedings of the IEEE 28th International Conference on Tools with Artificial Intelligence, San Jose, CA, USA, 6–8 November 2016; pp. 362–369.
- Yao, Y.; Li, X.; Liu, X.; Liu, P.; Liang, Z.; Zhang, J.; Mai, K. Sensing spatial distribution of urban land use by integrating points-of-interest and Google Word2Vec model. *Int. J. Geogr. Inf. Sci.* **2017**, *31*, 825–848. [\[CrossRef\]](#)
- Song, C.; Qu, Z.; Blumm, N.; Barabási, A.L. Limits of predictability in urban mobility. *Science* **2010**, *327*, 1018–1021. [\[CrossRef\]](#)
- Liu, Y.; Liu, X.; Gao, S.; Gong, L.; Kang, C.; Zhi, Y.; Shi, L. Social Sensing: A new approach to Understanding Our Socioeconomic Environments. *Ann. Assoc. Am. Geogr.* **2015**, *105*, 512–530. [\[CrossRef\]](#)
- Yin, J.; Dong, J.; Hamm, N.; Li, Z.; Wang, J.; Xing, H.; Fu, P. Integrating remote sensing and geospatial big data for urban land use mapping: A review. *Int. J. Appl. Earth. Obs.* **2021**, *103*, 102514. [\[CrossRef\]](#)
- Yuan, J.; Zheng, Y.; Xie, X. Discovering regions of different functions in a city using human mobility and POIs. In Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Beijing, China, 12–16 August 2012; pp. 186–194.
- Song, C.; Pei, T.; Ma, T.; Du, Y.; Shu, H.; Guo, S.; Fan, Z. Detecting arbitrarily shaped clusters in origin-destination flows using ant colony optimization. *Int. J. Geogr. Inf. Sci.* **2019**, *33*, 134–154. [\[CrossRef\]](#)
- Zhang, X.; Xu, Y.; Tu, W.; Ratti, C. Do different datasets tell the same story about urban mobility—A comparative study of public transit and taxi usage. *J. Transp. Geogr.* **2018**, *70*, 78–90. [\[CrossRef\]](#)

12. Zhai, W.; Bai, X.; Shi, Y.; Han, Y.; Peng, Z.R.; Gu, C. Beyond Word2vec: An approach for urban functional region extraction and identification by combining Place2vec and POIs. *Comput. Environ. Urban Syst.* **2019**, *74*, 1–12. [\[CrossRef\]](#)
13. Hu, S.; He, Z.; Wu, L.; Yin, L.; Xu, Y.; Cui, H. A framework for extracting urban functional regions based on multiprototype word embeddings using points-of-interest data. *Comput. Environ. Urban Syst.* **2020**, *80*, 101442. [\[CrossRef\]](#)
14. Ye, C.; Zhang, F.; Mu, L.; Gao, Y.; Liu, Y. Urban function recognition by integrating social media and street-level imagery. *Environ. Plan. B-Urban Anal. City Sci.* **2021**, *48*, 1430–1444. [\[CrossRef\]](#)
15. Yue, M.; Kang, C.; Andris, C.; Qin, K.; Liu, Y.; Meng, Q. Understanding the interplay between bus, metro, and cab ridership dynamics in Shenzhen, China. *Trans. GIS* **2018**, *22*, 855–871. [\[CrossRef\]](#)
16. Tu, W.; Zhu, T.; Xia, J.; Zhou, Y.; Lai, Y.; Jiang, J.; Li, Q. Portraying the spatial dynamics of urban vibrancy using multi-source urban big data. *Comput. Environ. Urban Syst.* **2020**, *80*, 101428. [\[CrossRef\]](#)
17. Liu, J.; Li, J.; Li, W.; Wu, J. Rethinking big data: A review on the data quality and usage issues. *ISPRS-J. Photogramm. Remote Sens.* **2016**, *115*, 134–142. [\[CrossRef\]](#)
18. Zhang, C.; Fu, H.; Hu, Q.; Cao, X.; Xie, Y.; Tao, D.; Xu, D. Generalized Latent Multi-View Subspace Clustering. *IEEE Trans. Pattern Anal. Mach. Intell.* **2020**, *42*, 86–99. [\[CrossRef\]](#) [\[PubMed\]](#)
19. Liu, Q.; Huan, W.; Deng, M.; Zheng, X.; Yuan, H. Inferring Urban Land Use from Multi-Source Urban Mobility Data Using Latent Multi-View Subspace Clustering. *ISPRS Int. J. Geo-Inf.* **2021**, *10*, 274. [\[CrossRef\]](#)
20. Sagioglu, S.; Sinanc, D. Big Data: A Review. In Proceedings of the 2013 International Conference on Collaboration Technologies and Systems, San Diego, CA, USA, 20–24 May 2013; pp. 42–47.
21. Fan, Y.; He, R.; Hu, B.G. Global and local consistent multi-view subspace clustering. In Proceedings of the Asian Conference on Pattern Recognition (ACPR), Kuala Lumpur, Malaysia, 3–6 November 2015; pp. 564–568.
22. Zhou, T.; Zhang, C.; Peng, X.; Bhaskar, H.; Yang, J. Dual Shared-Specific Multi-view Subspace Clustering. *IEEE T. Cybern.* **2019**, *50*, 3517–3530. [\[CrossRef\]](#)
23. Zheng, Q.; Zhu, J.; Ma, Y.; Li, Z.; Tian, Z. Multi-view subspace clustering networks with local and global graph information. *Neurocomputing* **2021**, *449*, 15–23. [\[CrossRef\]](#)
24. Tschannen, M.; Bachem, O.; Lucic, M. Recent advances in autoencoder-based representation learning. *arXiv* **2018**, arXiv:1812.05069.
25. Jarvis, R.A.; Patrick, E.A. Clustering using a similarity measure based on shared near neighbors. *IEEE Trans. Comput.* **1973**, *100*, 1025–1034. [\[CrossRef\]](#)
26. Toole, J.L.; Ulm, M.; González, M.C.; Bauer, D. Inferring land use from mobile phone activity. In Proceedings of the ACM SIGKDD International Workshop on Urban Computing, Beijing, China, 12–16 August 2012; pp. 1–8.
27. Krishna, K.; Murty, M. Genetic K-means algorithm. *IEEE Trans. Syst. Man Cybern.* **1999**, *29*, 433–439. [\[CrossRef\]](#) [\[PubMed\]](#)
28. Ng, A.; Jordan, M.; Weiss, Y. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems*; MIT Press: Vancouver, BC, Canada, 2002; Volume 14, pp. 849–856.
29. Ester, M.; Krieger, H.P.; Sander, J.; Xu, X. A density-based algorithm for discovering clusters in large spatial databases with noise. In Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining; 1996; Volume 96, pp. 226–231.
30. Yuan, N.J.; Zheng, Y.; Xie, X.; Wang, Y.; Zheng, K.; Xiong, H. Discovering urban functional zones using latent activity trajectories. *IEEE Trans. Knowl. Data Eng.* **2015**, *27*, 712–725. [\[CrossRef\]](#)
31. Gao, H.; Nie, F.; Li, X.; Huang, H. Multi-view subspace clustering. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), Santiago, Chile, 13–16 December 2015; pp. 4238–4246.
32. Parsons, L.; Haque, E.; Liu, H. Subspace clustering for high dimensional data: A review. *Acm Sigkdd Explor. Newsl.* **2004**, *6*, 90–105. [\[CrossRef\]](#)
33. Vidal, R. Subspace clustering. *IEEE Signal. Process. Mag.* **2011**, *28*, 52–68. [\[CrossRef\]](#)
34. Liu, G.; Lin, Z.; Yan, S.; Sun, J.; Yu, Y.; Ma, Y. Robust recovery of subspace structures by low-rank representation. *IEEE Trans. Pattern Anal. Mach. Intell.* **2012**, *35*, 171–184. [\[CrossRef\]](#)
35. Elhamifar, E.; Vidal, R. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE Trans. Pattern Anal. Mach. Intell.* **2013**, *35*, 2765–2781. [\[CrossRef\]](#)
36. Hu, H.; Lin, Z.; Feng, J.; Zhou, J. Smooth representation clustering. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Columbus, OH, USA, 24–17 June 2014; pp. 3834–3841.
37. Li, C. Structured sparse subspace clustering: A joint affinity learning and subspace clustering framework. *IEEE Trans. Image Process.* **2017**, *26*, 2988–3001. [\[CrossRef\]](#)
38. Cao, X.; Zhang, C.; Fu, H.; Liu, S.; Zhang, H. Diversity-induced multi-view subspace clustering. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Boston, MA, USA, 7–12 June 2015; pp. 586–594.
39. Luo, S.; Zhang, C.; Zhang, W.; Cao, X. Consistent and specific multi-view subspace clustering. In Proceedings of the AAAI Conference on Artificial Intelligence, New Orleans, LA, USA, 2–7 February 2018; pp. 3730–3737.
40. Zhu, P.; Hui, B.; Zhang, C.; Du, D.; Wen, L.; Hu, Q. Multi-view Deep Subspace Clustering Networks. *arXiv* **2019**, arXiv:1908.01978. 2019.
41. Zhang, C.; Hu, Q.; Fu, H.; Zhu, P.; Cao, X. Latent multi-view subspace clustering. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 4279–4287.
42. Yu, X.; Liu, H.; Wu, Y.; Zhang, C. Intrinsic self-representation for multi-view subspace clustering. *Sci. China Inf. Sci.* **2021**, *51*, 1625–1639.

43. Wang, X.; Liu, H.; Qian, X.; Jiang, Y.; Deng, Z.; Wang, S. Cascaded hidden space feature mapping, fuzzy clustering, and nonlinear switching regression on large datasets. *IEEE Trans. Fuzzy Syst.* **2018**, *26*, 640–655. [\[CrossRef\]](#)
44. Wang, X.; Lei, Z.; Guo, X.; Zhang, C.; Shi, H.; Li, S.Z. Multi-view subspace clustering with intactness-aware similarity. *Pattern Recognit.* **2019**, *88*, 50–63. [\[CrossRef\]](#)
45. Zhu, W.; Lu, J.; Zhou, J. Structured General and Specific Multi-view Subspace Clustering. *Pattern Recognit.* **2019**, *93*, 392–403. [\[CrossRef\]](#)
46. Zheng, Q.; Zhu, J.; Li, Z.; Pang, S.; Wang, J.; Li, Y. Feature concatenation multi-view subspace clustering. *Neurocomputing* **2020**, *379*, 89–102. [\[CrossRef\]](#)
47. Xia, S.; Xiong, Z.; Luo, Y.; Zhang, G. Effectiveness of the Euclidean distance in high dimensional spaces. *Optik* **2015**, *126*, 5614–5619. [\[CrossRef\]](#)
48. Liu, Q.; Liu, W.; Deng, M.; Cai, J.; Liu, Y. An adaptive detection of multilevel co-location patterns based on natural neighborhoods. *Int. J. Geogr. Inf. Sci.* **2021**, *35*, 556–581. [\[CrossRef\]](#)
49. Wang, Q.; Cheng, J.; Gao, Q.; Zhao, G.; Jiao, L. Deep multi-view subspace clustering with unified and discriminative learning. *IEEE Trans. Multimed.* **2020**, *23*, 3483–3493. [\[CrossRef\]](#)
50. Comito, C.; Talia, D. GDIS: A service-based architecture for data integration on Grids. In *On the Move to Meaningful Internet Systems 2004: OTM 2004 Workshops*; Meersman, R., Tari, Z., Corsaro, A., Eds.; OTM 2004. Lecture Notes in Computer Science; Springer: Berlin, Heidelberg, 2004; Volume 3292.
51. Lee, J.; Kang, M. Geospatial big data: Challenges and opportunities. *Big Data Res.* **2015**, *2*, 74–81. [\[CrossRef\]](#)
52. Liu, X.; Tian, Y.; Zhang, X.; Wan, Z. Identification of urban functional regions in chengdu based on taxi trajectory time series data. *ISPRS Int. J. Geo-Inf.* **2020**, *9*, 158. [\[CrossRef\]](#)
53. Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G.S.; Dean, J. Distributed representations of words and phrases and their compositionality. In *Proceedings of the Advances in Neural Information Processing Systems*, Lake Tahoe, CA, USA, 9–12 December 2013; p. 26.
54. Lau, J.; Baldwin, T. An empirical evaluation of doc2vec with practical insights into document embedding generation. *arXiv* **2016**, arXiv:1607.05368.
55. Ertöz, L.; Steinbach, M.; Kumar, V. Finding clusters of different sizes, shapes, and densities in noisy, high dimensional data. In *Proceedings of the 2003 SIAM International Conference on Data mining*, Society for Industrial and Applied Mathematics, San Francisco, CA, USA, 1–3 May 2003; pp. 47–58.
56. Tan, P.N.; Steinbach, M.; Karpatne, A.; Kumar, V. *Introduction to Data Mining*; Pearson Education: London, UK, 2006; pp. 622–630.
57. Liu, Q.; Deng, M.; Bi, J.; Yang, W. A novel method for discovering spatio-temporal clusters of different sizes, shapes, and densities in the presence of noise. *Int. J. Digit. Earth* **2014**, *7*, 138–157. [\[CrossRef\]](#)
58. Kingma, D.; Ba, J. Adam: A method for stochastic optimization. *arXiv* **2014**, arXiv:1412.6980.
59. Rousseeuw, P.J. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.* **1987**, *20*, 53–65. [\[CrossRef\]](#)
60. Martínez, L.; Viegas, J.; Silva, E. A traffic analysis zone definition: A new methodology and algorithm. *Transportation* **2009**, *36*, 581–599. [\[CrossRef\]](#)
61. Yang, B.; Tian, Y.; Wang, J.; Hu, X.; An, S. How to improve urban transportation planning in big data era? A practice in the study of traffic analysis zone delineation. *Transp. Policy* **2022**, *127*, 1–14. [\[CrossRef\]](#)
62. Cherry, J.M.; Adler, C.; Ball, C.; Chervitz, S.A.; Dwight, S.S.; Hester, E.T.; Botstein, D. SGD: Saccharomyces genome database. *Nucleic Acids Res.* **1998**, *26*, 73–79. [\[CrossRef\]](#)
63. Xia, R.; Pan, Y.; Du, L.; Yin, J. Robust multi-view spectral clustering via low-rank and sparse decomposition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Québec, QC, Canada, 27–31 July 2014; Volume 28.