

1 Article

2 An Efficient Grid-based K-prototypes Algorithm for 3 Sustainable Decision Making using Spatial Objects

4 Hong-Jun Jang ¹, Byoungwook Kim ², Jongwan Kim ³ and Soon-Young Jung ^{4,*}

5 ¹ Department of Computer Science and Engineering, Korea University, Seoul, 02841, Korea;

6 hongjunjang@korea.ac.kr

7 ² Department of Computer Engineering, Dongguk University, Gyeongju, 38066, Korea;

8 bwkim@dongguk.ac.kr

9 ³ Smith Liberal Arts College, Sahmyook University, Seoul, 01795, Korea; kimj@syu.ac.kr

10 ⁴ Department of Computer Science and Engineering, Korea University, Seoul, 02841, Korea; jsy@korea.ac.kr

11 * Correspondence: jsy@korea.ac.kr; Tel.: +82-2-3290-2394

12

13 **Abstract:** Data mining plays a critical role in the sustainable decision making. The k-prototypes
14 algorithm is one of the best-known algorithm for clustering both numeric and categorical data.
15 Despite this, however, clustering a large number of spatial object with mixed numeric and
16 categorical attributes is still inefficient due to its high time complexity. In this paper, we propose an
17 efficient grid-based k-prototypes algorithms, GK-prototypes, which achieves high performance for
18 clustering spatial objects. The first proposed algorithm utilizes both maximum and minimum
19 distance between cluster centers and a cell, which can remove unnecessary distance calculation. The
20 second proposed algorithm as extensions of the first proposed algorithm utilizes spatial dependence
21 that spatial data tend to be more similar as objects are closer. Each cell has a bitmap index which
22 stores categorical values of all objects in the same cell for each attribute. This bitmap index can
23 improve the performance in case that a categorical data is skewed. Our evaluation experiments
24 showed that proposed algorithms can achieve better performance than the existing pruning
25 technique in the k-prototypes algorithm.

26 **Keywords:** clustering; spatial data; grid-based k-prototypes; data mining; sustainability

27

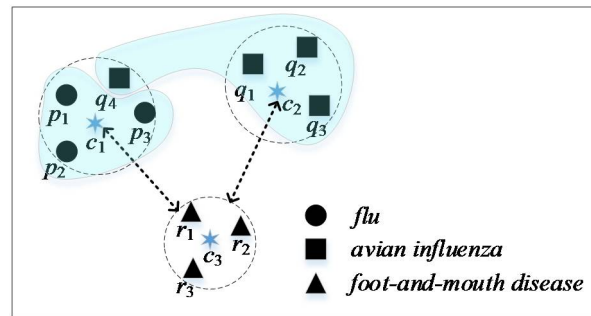
28 1. Introduction

29 Sustainability is a concept for balancing environmental, economic and social dimensions with
30 decision-making [23]. Data mining in sustainability is a very important issue since it sustainable
31 decision making contributes to the transition to a sustainable society [24]. Especially, there is a
32 growing interest in spatial data mining in making sustainable decisions in geographical
33 environments and national land policies [25].

34 Recently, spatial data mining has become more and more important as spatial data collection is
35 increasing due to technological developments such as geographic information system (GIS) and
36 global positioning system (GPS) [26][27]. The main techniques of spatial data mining are spatial
37 clustering [1], spatial classification [2], spatial association rule [3], and spatial characterization [4].
38 Spatial clustering is a technique used to classify data with high similar geographic and locational
39 characteristics into the same group. It is an important component to discover hidden knowledge in a
40 huge of spatial data [5]. Spatial clustering is used in the HotSpot detection which detects areas where
41 specific events occur [6]. Hotspot detection is used in various fields such as crime analysis [7,8,9], fire
42 analysis [10] and disease analysis [11,12,13,14].

43 Most spatial clustering studies have focused on efficiently finding groups for numeric data such
44 as location information of spatial objects. However, many real-world spatial objects have categorical
45 data, not just numeric data. Hence, if the categorical data affect spatial clustering results, the error
46 value can be increased when the final cluster results are evaluated.

47



48

49

Figure 1. An example of clustering using location information

50

51

52

53

54

55

56

57

We present an example of disease analysis clustered using only location information. Figure 1 shows 10 cases of disease location divided into three clusters. Diseases are divided into three categories (i.e. p , q and r), and the measures to be taken in the area vary according to each disease. In Figure 1, the representative attribute of cluster c_i is set to p , so we only deal with p . Therefore, when q occurs in the same area, it is difficult to cope with. In this example, three clusters are constructed using only numeric data (location information of disease occurrence). If the data from this example was used to inform a policy decision, it could result in a decision maker failing to implement the correct policy.

58

59

60

61

62

63

64

65

66

67

Data generated in the real world is often mixed with numeric data as well as categorical data. In order to apply the clustering technique to the real world, algorithms that can consider categorical data are required. A representative clustering algorithm that can use mixed data is the k-prototypes algorithm [15]. The basic k-prototypes algorithm has a large time complexity due to the processing of all data. Therefore, it is important to reduce execution time in order to process the k-prototypes algorithm on large data. However, only a few studies have been conducted to reduce the time complexity of the k-prototypes algorithm. Kim [16] proposed a pruning technique to reduce distance computation between an object and cluster centers using the concept of partial distance computation. However, this method does not have high pruning efficiency by comparing objects one by one with cluster center.

68

69

70

71

To improve performance, we propose an effective grid-based k-prototypes algorithm, GK-prototypes, for clustering spatial objects. The proposed method makes use of the grid-based indexing technique which improve pruning efficiency to compare distance between cluster centers and a cell instead of cluster centers and an object.

72

73

74

75

76

77

78

79

Spatial data can have geographic data as categorical attributes that indicate the characteristics of the object as well as the location of the object. Geographic data tend to have spatial dependence. Spatial dependence is the property of objects that are close to each other having increased similarities [17]. For example, soil types or network type are more likely to be similar at points one meter apart than at points one kilometer apart. Due to the nature of spatial dependence, the categorical data of spatial data is often skewed according to the position of the object. For improving performance of a grid-based k-prototypes algorithm, we take advantage of the spatial dependence to the bitmap indexing technique.

80

The contributions of this paper are summarized as follows.

81

82

83

84

85

86

87

88

89

- We proposed an effective grid-based k-prototypes, GK-prototypes, which improve the performance of a basic k-prototypes algorithm.
- We developed a pruning technique which utilizes the minimum and maximum distance on numeric attributes and the maximum distance on categorical attributes between a cell and a cluster center.
- We developed a pruning technique based on a bitmap index to improve the efficiency of the pruning in case that a categorical data is skewed.
- We conducted several experiments on synthetic datasets. Our algorithms can achieve better performance than the existing pruning technique in the k-prototypes algorithm.

90 The organization of the rest of this paper is as follows. In Section 2, the basic k-prototypes
 91 algorithm and the previous research on pruning in the k-prototypes algorithm are described. In
 92 Section 3, we first briefly describe some basic notations and definitions. After that, the proposed GK-
 93 prototypes algorithm is explained in Section 4. In Section 5, experimental results on synthetic data
 94 demonstrate the performance. Section 6 concludes the paper.

95 2. Related works

96 2.1. The k-prototypes algorithm

97 The k-prototypes algorithm is first proposed clustering algorithm to deal with mixed data types
 98 (numeric data and categorical data), which integrates k-means and k-modes algorithms [15]. Let a set
 99 of n objects be $O=\{o_1, o_2, \dots, o_n\}$ where $o_i=(o_{i1}, o_{i2}, \dots, o_{im})$ is consisted of m attributes. The purpose of
 100 clustering is to partition n objects into k disjoint clusters $C=\{C_1, C_2, \dots, C_n\}$ according to the degree
 101 of similarity of objects. The distance is used as a measure to group objects with high similarity into the
 102 same cluster. The distance $d(o_i, C_j)$ between o_i and C_j is calculated as follows:

$$d(o_i, C_j) = d_r(o_i, C_j) + d_c(o_i, C_j) \quad (1)$$

103 where $d_r(o_i, C_j)$ is the distance between numeric attributes and $d_c(o_i, C_j)$ is the distance
 104 between categorical attributes.

$$d_r(o_i, C_j) = \sum_{k=1}^p |o_{ik} - c_{jk}|^2 \quad (2)$$

$$d_c(o_i, C_j) = \sum_{k=p+1}^m \delta(o_{ik}, c_{jk}) \quad (3)$$

$$\delta(o_{ik}, c_{jk}) = \begin{cases} 0, & \text{when } o_{ik} = c_{jk} \\ 1, & \text{when } o_{ik} \neq c_{jk} \end{cases} \quad (4)$$

105 In Equation (2), $d_r(o_i, C_j)$ is the squared Euclidean distance between an object and a cluster center
 106 on the numeric attributes. $d_c(o_i, C_j)$ is the dissimilar distance on the categorical attributes, where o_{ik}
 107 and c_{jk} , $1 \leq k \leq p$, are values of numeric attributes, o_{ik} and c_{jk} , $p+1 \leq k \leq m$ are values of categorical attributes.
 108 That is, p is the number of numeric attributes and $m-p$ is the number of categorical attributes.

109 2.2. Existing pruning technique in the k-prototypes algorithm

110 The k-prototypes algorithm spends most of execution time computing the distance between an
 111 object and cluster centers. In order to improve the performance of the k-prototypes algorithm, Kim
 112 [16] proposed the concept of partial distance computation (PDC) which compares only partial
 113 attributes, not all attributes in measuring distance. The maximum distance that can be measured in
 114 one categorical attribute is 1. Thus the distance that can be measured from the categorical attributes
 115 is bound to the number of categorical attributes, $m-p$. Given an object o and the two cluster centers (c_1
 116 and c_2), if the difference between $d_r(o, c_1)$ and $d_r(o, c_2)$ is more than $m-p$, we can know which clusters
 117 are closer to the object without the distance using the categorical attributes. However, PDC is still not
 118 efficient due to the fact that all objects are involved in the distance calculation and the characteristic
 119 (i.e. spatial dependence) of spatial data is not utilized in the clustering process.

120 2.3. Grid-based clustering algorithm

121 The grid-based techniques have the fastest processing time that depends on the number of the
 122 grid cells instead of the number of objects in the data set [18]. The basic grid-based algorithm is as
 123 follows. At first, a set of grid-cells is defined. In general, these grid-based clustering algorithm use a
 124 single uniform or multi-resolution grid cell to partition the entire datasets into cells. Each object is
 125 assigned to the appropriate grid cell and the density of each cell is computed. The cells, whose a

126 degree of density is below a certain threshold, are eliminated. In addition to the density of cells,
 127 statistical information of objects in the cell is computed. After that, the clustering process is performed
 128 on the grid cells using each cell's statistical information, instead of the objects itself.

129 The representative grid-based clustering algorithms are STING [19] and CLIQUE [20]. STING is
 130 a grid-based multi resolution clustering algorithm in which the spatial area is divided into
 131 rectangular cells with a hierarchical structure. Each cell at a high level is divided into several smaller
 132 cells in the next lower level. For each cell in pre-selected layer, the relevancy of the cell is checked by
 133 computing the confidence interval. If the cell is relevant, we include the cell in a cluster. If the cell is
 134 irrelevant, it is removed from further consideration. We look for relevant cells at the next lower layer.
 135 This algorithm combines relevant cells into relevant regions and return the so obtained clusters.
 136 CLIQUE is a grid-based and density-based clustering algorithm to identify subspaces of a high
 137 dimensional data that allow better clustering quality than original data. CLIQUE partitions the n -
 138 dimensional data into non overlapping rectangular units. The units are obtained by partitioning
 139 every dimension into certain intervals of equal length and selectivity of a unit is defined as the total
 140 data points contained in it. A cluster in CLIQUE is a maximal set of connected dense units within a
 141 subspace. In a grid-based clustering study, the grid is used in order that clustering is performed on
 142 the grid cells, instead of objects itself. Chen et al. [21] proposes algorithm called GK-means, which
 143 integrates grid structure and spatial index with k-means algorithm. It focuses on choice the better
 144 initial centers to improve the clustering quality and to reduce the computational complexity of k-
 145 means.

146 Most existing grid-based clustering algorithms regard objects in same cell of grid as a data point
 147 to process large scale data. Thus, the final clustering results of these algorithms are not the same as a
 148 basic k-prototype cluster result, but all the cluster boundaries are either horizontal or vertical. In GK-
 149 means, the grid is used to select initial centers and remove noise data, but not used to reduce
 150 unnecessary distance calculation. To the best of our knowledge, such a grid-based pruning technique
 151 to improve the performance of the k-prototypes algorithm has not been previously demonstrated.

152 3. Preliminary

153 In this section, we present some basic notations and definitions before describing our algorithms.
 154 We summarize the notation used throughout this paper in Table 1.

155 **Table 1.** A summary of notations

Notation	Description
O	a set of data
o_i	i -th data in O
n	the number of objects
m	the number of attributes of an object
c_i	the i cluster center point
v_i	the value of grid partition interval
g^k	a cell of grid
$d(o_i, c_j)$	a distance between an object and an cluster center
$d_r(o_i, c_j)$	a distance between an object and an cluster center for only numeric attributes
$d_c(o_i, c_j)$	a distance between an object and an cluster center for only categorical attributes
$d_{min}(g_i, c_j)$	the minimum distance between a cell and a cluster center for only numeric attributes
$d_{max}(g_i, c_j)$	the maximum distance between a cell and a cluster center for only numeric attributes

156 Consider a set of n objects, $O=\{o_1, o_2, \dots, o_n\}$. $o_i=(o_{i1}, o_{i2}, \dots, o_{im})$ is an object represented by m
 157 attribute values. The m attributes consist of m_r (the number of numeric attributes) and m_c (the number
 158 of categorical attributes), $m=m_r+m_c$. The distance between an object and a cluster center, $d(o, c)$, is
 159

160 calculated by Equation (1) in Section 2. We adopt the data indexing technique based on grid for
 161 pruning. First we define the cells that make up the grid.

162
 163 **Definition 1.** A cell g in m_r -dimension grid is defined by a start point vector S and an end point vector
 164 T : $g = (S, T)$, where $S = [s_1, s_2, \dots, s_{m_r}]$ and $T = [t_1, t_2, \dots, t_{m_r}]$ and $s_i \leq t_i$ for $1 \leq i \leq m_r$ and $s_i + v_i = t_i$. The v_i
 165 is interval distance between start and end position of a cell g on i -dimension.
 166

167 **Definition 2.** The minimum distance between a cell g_i and a cluster center c_j for numeric attributes,
 168 denoted $d_{min}(g_i, c_j)$, is;

$$169 \quad d_{min}(g_i, c_j) = \sqrt{\sum_{i=1}^{m_r} |o_i - r_i|^2},$$

$$170 \quad \text{where } r_i = \begin{cases} s_i & \text{if } o_i < s_i \\ t_i & \text{if } o_i > t_i \\ o_i & \text{otherwise.} \end{cases}$$

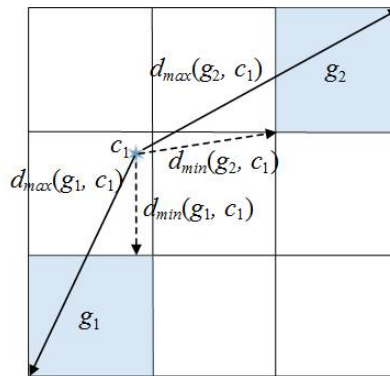
171 We use the classic Euclidian distance to measure the distance. If a cluster center is inside the cell,
 172 the distance between them is zero. If a cluster center is outside the cell, we use the Euclidean distance
 173 between the cluster center and the nearest edge of the cell.
 174

175 **Definition 3.** The maximum distance between a cell g_i and a cluster center c_j for numeric attributes,
 176 denoted $d_{max}(g_i, c_j)$, is;

$$177 \quad d_{max}(g_i, c_j) = \sqrt{\sum_{i=1}^{m_r} |p_i - r_i|^2},$$

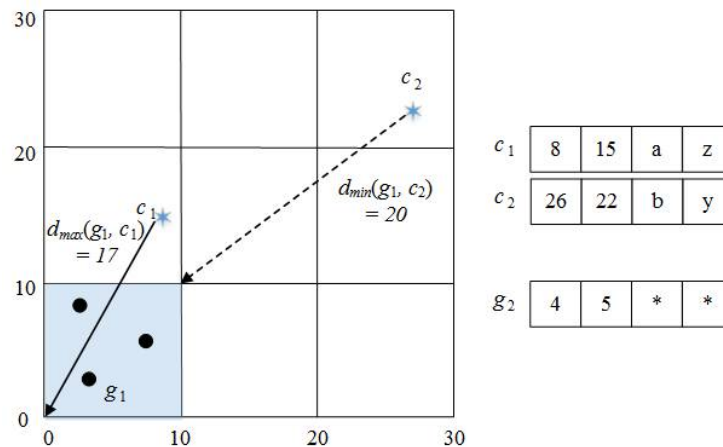
$$178 \quad \text{where } r_i = \begin{cases} t_i, & p_i \leq \frac{s_i + t_i}{2} \\ s_i, & \text{otherwise.} \end{cases}$$

179 To distinguish between the two distances d_{min} and d_{max} , an example is illustrated in Figure 2,
 180 showing a cluster center (c_1), two cells (g_1 and g_2) and the corresponding distances.
 181



182
 183 **Figure 2.** An example of distances d_{min} and d_{max} .

184 In a grid-based approach, $d_{min}(g,c)$ and $d_{max}(g,c)$ for a cell g and cluster centers c , are measured
 185 firstly before measuring the distance between an object and cluster centers. We can use d_{min} and d_{max}
 186 to improve performance of k-prototypes algorithm. Figure 3 shows an example of pruning using d_{min}
 187 and d_{max} . The object consists of 4 attributes (2 numeric and 2 categorical attributes).



188

189

Figure 3. An example of pruning method using minimum distance and maximum distance

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

In Figure 3, the maximum distance between g_1 and c_1 , $d_{max}(g_1, c_1) = 17(\sqrt{(8-0)^2 + (15-0)^2})$. The minimum distance between g_1 and c_2 , $d_{min}(g_1, c_2) = 20(\sqrt{(26-10)^2 + (22-10)^2})$. The $d_{max}(g_1, c_1)$ is three less than $d_{min}(g_1, c_2)$. Therefore, all objects in g_1 are closer to c_1 than c_2 , if we are considering only numeric attributes. To find the closest cluster center from an object, we have to measure the distance by categorical attributes. If the difference between d_{min} and d_{max} is more than m_c (maximum distance by categorical attributes), however, the cluster closest to the object can be determined without the distance by categorical attributes. In Figure 3, the categorical data of c_1 is (a, z), and the categorical data of c_2 is (b, y). Assume that there are no objects in g_1 with 'a' in the first categorical attribute and 'z' in second categorical attribute. Since $d_c(o, c_1)$ of all objects in g_1 is 2, maximum distance of all objects in g_1 is $d_{max}(g_1, c_1) + 2$. The maximum distance between c_1 and objects in g_1 is not less than the minimum distance between c_2 and objects in g_1 . We can know that all objects in g_1 are closer to c_1 than c_2 .

Lemma 1. For any cluster centers c_i, c_j and any cell g_x , if $d_{min}(g_x, c_j) - d_{max}(g_x, c_i) > m_c$ then, $\forall o \in g_x, d(o, c_i) < d(o, c_j)$.

Proof. By assumption, $d_{min}(g_x, c_j) > m_c + d_{max}(g_x, c_i)$. By Definition 2 and 3, $d_{min}(g_x, c_i) \leq d(c_i, o) \leq d_{max}(g_x, c_i) + m_c$, and $d_{min}(g_x, c_j) \leq d(c_j, o) \leq d_{max}(g_x, c_j) + m_c$. $d(c_i, o) \leq d_{max}(g_x, c_i) + m_c < d_{min}(g_x, c_j) \leq d(c_j, o)$. $\therefore \forall o \in g_x, d(c_i, o) < d(c_j, o)$. \square .

Lemma 1 is the basis for our proposed pruning techniques. In the process of clustering, we first exploit Lemma 1 to remove cluster centers to be compared to objects.

209

4. GK-prototypes algorithm

210

211

212

213

214

In this section, we present two pruning techniques that are based on grid for improving the performance of the k-prototypes algorithm. The first pruning technique is KCP (K-prototypes algorithm with Cell Pruning) which utilizes d_{min} , d_{max} and the maximum distance on categorical attributes. The second pruning technique is KBP (K-prototypes algorithm with Bitmap Pruning) which utilizes bitmap indexes to reduce unnecessary distance calculation on categorical attributes.

215

4.1. Cell pruning technique

216

217

218

219

220

221

222

223

The computational cost of the k-prototypes algorithm is most often encountered in the step of measuring distance between objects and cluster centers. To improve the performance, each object is indexed into a grid by numeric data in data preparation step. We set up grid cells storing two types of information. a) The first is a start point vector S and an end point vector T , which is the range of the numeric value of the objects to be included in the cell (see Definition 1). Based on this cell information, the minimum and maximum distances between each cluster centers and a cell are measured. b) The second is bitmap indexes which is explained in Subsection 4.2.

224 **Algorithm 1** The k-prototypes algorithm with cell pruning (KCP)
 225 Input: k: the number of cluster, G: the grid in which all objects are stored per cell
 226 Output: k cluster centers

```

227 1: C[ ] ← ∅ // k cluster centers
228 2: Randomly choosing k object, and assigning it to C.
229 3: while IsConverged() do
230 4:   for each cell g in G
231 5:     dmin[ ], dmax[ ] ← Calc(g, C)
232 6:     dminmax ← min(dmax[ ])
233 7:     candidate ← ∅
234 8:     for j ← 1 to k do
235 9:       if (dminmax + mc > dmin[j]) // Lemma 1
236 10:        candidate ← candidate ∪ j
237 11:     end for
238 12:     min_distance ← ∞
239 13:     min_cluster ← null
240 14:     for each object o in g
241 15:       for each center c in candidate
242 16:         if min_distance > d(o, c)
243 17:           min_cluster ← index of c
244 18:       end for
245 19:       Assign(o, min_cluster)
246 20:       UpdateCenter(C[c])
247 21:     end for
248 22: end while
249 23: return C[k]

```

250
 251 The details of the cell pruning are as follows. First, we initialize an array C[j] to store the position
 252 of k cluster centers, $1 \leq j \leq k$. Various initial cluster centers selection methods have been studied to
 253 improve the accuracy of clustering results in the k-prototypes algorithm [22]. Since we aim at
 254 increasing the clustering speed improvement, however, we adopt a simple method to select k objects
 255 randomly from input data and use them as initial cluster centers.

256 In general, the result of clustering algorithm is evaluated after a single clustering process has
 257 been performed. Based on these evaluation result, it is determined whether the same clustering
 258 process have to be repeated or terminated. The iteration is terminated if the sum of the difference
 259 between the current cluster center and the previous cluster center is less than a predetermined value
 260 (ϵ) as an input parameter by users. In Algorithm 1, we determine the termination condition through
 261 the *IsConverged()* function of the while statement.

262 In iteration step of clustering (line 4), the distances between objects and k cluster centers are
 263 measured by cells. The *Calc(g, C)* function returns the minimum and maximum distances between a
 264 cell g and k cluster centers (C) for each cluster center (line 5). The smallest distance among the
 265 maximum distance is stored in the *dminmax* (line 6). The candidate stores the index number of the
 266 cluster center that need to be measured from the cell. Through Lemma 1, if *dminmax*+*m_c* is greater
 267 than *dmin[j]*, the j cluster center is included in the distance calculation, otherwise it is excluded (lines
 268 8-11). After Lemma 1 is applied, only the cluster centers to be measured the distance from objects in
 269 the cell are finally left in the candidate. All objects in the cell are measured from the cluster centers in
 270 the candidate, *d(o, c)*, (line 16). An object is assigned to the cluster where *d(o, c)* is computed as the
 271 smallest value using *Assign(o, min_cluster)* function. The center of cluster to which a new object is
 272 added is updated by remeasuring its cluster center using *UpdateCenter(C[c])* function. This clustering
 273 process is repeated until the end condition *IsConverged()* is return a false.

274 4.2. Bitmap pruning technique

275 Spatial data tend to have similar categorical values in neighboring objects, categorical data is
 276 often skewed. If we can utilize the characteristics of spatial data, the performance of the k-prototypes
 277 algorithm can be further improved. In this subsection, we introduce the KBP that can improve the
 278 efficiency of pruning when categorical data is skewed.

279 The KBP stores categorical data of all objects in a cell in a bitmap index. Figure 4 shows an
 280 example of storing categorical data as a bitmap index. Figure 4(a) is an example of spatial data. The
 281 x and y attributes indicate location information of objects as numeric data, and z and w attributes
 282 indicate features of objects as categorical data. For five objects in the same cell, $g=\{o_1, o_2, o_3, o_4, o_5\}$,
 283 Figure 4(b) shows the bitmap index structure where a row presents a categorical attribute and a
 284 column presents a categorical data of objects in same cell. A bitmap index consists of one vector of
 285 bits per attribute value, where the size of each bitmap is equal to the number of categorical data in
 286 the raw data. The bitmaps are encoded such that the i -th bit is set to 1 if the raw data has a value
 287 corresponding to the i -th column in the bitmap index, otherwise it is set to 0. For example, the value
 288 1 in z row and c column from bitmap index means that the value c exists in z attribute of raw data.
 289 When raw data is converted to a bitmap index, object id information is removed. We can quickly
 290 check for the existence of the specified value in raw data using the bitmap index.
 291

oid	x	y	z	w
o_1	14	6	A	D
o_2	12	3	F	B
o_3	30	4	C	D
o_4	31	9	F	C
o_5	59	9	F	E

(a) Example mixed data

	A	B	C	D	E	F
z	1	0	1	0	0	1
w	0	1	1	1	1	0

(b) Bitmap indexing structure

292

293

Figure 4. Bitmap indexing structure

294 A maximum of categorical distance is determined by the number of categorical attributes (m_c).
 295 If the difference of two numeric distance between one object and two cluster centers, $|dr(o, c_i) - dr(o,$
 296 $c_j)|$, is more than m_c , we can know the cluster center closer to the object without categorical distance.
 297 However, since the numeric distance is not known in advance, we cannot determine the cluster center
 298 closer to the object by only categorical distance. Thus, the proposed KBP is utilized in reducing
 299 categorical distance calculations in the KCP. Algorithm 2 describes the proposed KBP. We explain
 300 only the extended parts from the KCP.

301

302 **Algorithm 2** The k-prototypes algorithm with bitmap pruning (KBP)

303 Input: k: the number of cluster, G: the grid in which all objects are stored per cell

304 Output: k cluster centers

305 1: $C[k] \leftarrow \emptyset$ // k cluster center
 306 2: Randomly choosing k object, and assigning it to C.
 307 3: **while** $IsConverged()$ **do**
 308 4: **for each** cell g in G
 309 5: $dmin[], dmax[] \leftarrow Calc(g, C)$
 310 6: $dminmax \leftarrow min(dmax[])$
 311 7: $arrayContain[] \leftarrow IsContain(g, C)$
 312 8: $candidate \leftarrow \emptyset$
 313 9: **for** $j \leftarrow 1$ **to** k **do**
 314 10: **if** $(dminmax + m_c > dmin[j])$ // Lemma 1


```

315 11:     candidate ← candidate ∪ j
316 12:   end for
317 13:   min_distance ← ∞
318 14:   min_cluster ← null
319 15:   distance ← null
320 16:   for each object o in g
321 17:     for each center c in candidate
322 18:       if (arrayContain [c] == 0)
323 19:         distance =  $d_r(o, c) + m_c$ 
324 20:       else
325 21:         distance =  $d_r(o, c) + d_c(o, c)$ 
326 22:       if min_distance > distance
327 23:         min_cluster ← index of cluster center
328 24:     end for
329 25:     Assign(o, min_cluster)
330 26:     UpdateCenter(C)
331 27:   end for
332 28: end while
333 29: return C[k]

```

334

335 To improve the efficiency of pruning on categorical data, the KBP is implemented by adding
336 two step to the KCP. The first step is to find out whether the categorical attributes value of each cluster
337 centers exists in the bitmap index of the cell using the *IsContain* function (line 7). The *IsContain*
338 function compares the categorical data of each cluster center with the bitmap index and returns 1 if
339 there is more than one of the same data in the corresponding attribute. Otherwise, it returns 0. In
340 Figure 4 (b), assume that we have a cluster center, $c_i = (2, 5, D, A)$. The bitmap index does not have D
341 in the z attribute and A in the w attribute. In this case, we can know that there are no objects in the
342 cell that have D in the z attribute and A in the w attribute. Therefore, the maximum categorical
343 distance between all objects belonging to a cell and cluster center c_i is 0. Assume that we have another
344 the cluster center, $c_j = (2, 5, A, B)$. The bitmap index has A in the z attribute and B in the w attribute.
345 In this case, we can know that there are some objects in the cell that have A in the z attribute or b
346 in the w attribute. If one or more objects with the same categorical value are found in the corresponding
347 attribute, the categorical distance calculation has to be performed for all objects in the cell in order to
348 know the correct categorical distance. Finally, *arrayContain[i]* stores the result of the comparison
349 between the *i*-th cluster center and the cell. In lines 8-12, the cluster centers that need to measure
350 distance with each cell *g* are stored in the candidate like as the KCP. The second extended part (lines
351 18-23) is used to determine whether to measure the categorical distance using *arrayContain*. If
352 *arrayContain[i]* has 0, the m_c is directly used as the result of the categorical distance without measuring
353 the categorical distance (line 19).

354 5. Experiments

355 In this section, we evaluate the performance of the proposed GK-prototypes algorithms. For
356 performance evaluation, we compare partial distance computation pruning (PDC) [16] and our two
357 algorithms (KCP and KBP). We examined the performance of our algorithms as the number of objects,
358 the number of clusters, the number of categorical attributes and the number of division in each
359 dimension increased. In Table 2, the parameters used for the experiments are summarized. The
360 rightmost column of the table means the baseline values of the various parameters. For each set of
361 parameters, we perform 10 sets of experiments and the average values are reported. Even under the
362 same parameter values, the number of clustering iterations will vary if the input data is different.

363 Therefore, we measure the performance of each algorithm with the average execution time that it
 364 takes for clustering to repeat once.

365 **Table 2.** Parameters for the experiments

Parameter	Description	Baseline Value
n	no. of objects	1,000K
k	no. of clusters	10
m_r	no. of numeric attributes	2
m_c	no. of categorical attributes	5
s	no. of division in each dimensions	10

366
 367 The experiments are carried out on a PC with Intel(R) Core(TM) i7 3.5 GHz, 32GB RAM. All the
 368 algorithms are implemented in Java.

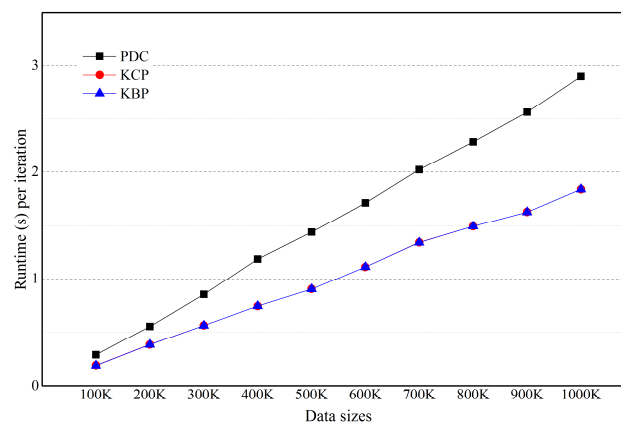
369 5.1. Data sets

370 We generate many synthetic datasets with numeric attributes and categorical attributes. For
 371 numeric attributes in each dataset, two numeric attributes are generated in the 2D space $[0, 100] \times [0,$
 372 $100]$ to indicate an object's location. Each object is assigned into $s \times s$ cells in a grid by these numeric
 373 data. The numeric data is generated according to uniform distributions in which each numeric data
 374 is selected in $[0, 100]$ randomly or Gaussian distributions with mean = 50 and standard deviation =
 375 10. The categorical data is generated according to uniform distributions or skewed distributions. For
 376 uniform distributions, we select an alphabet from A to Z randomly. For skewed distributions, we
 377 generate categorical data in such a way that objects in same cell have similar alphabet based on its
 378 numeric data.

379 5.2. Effects of the number of objects

380 To illustrate scalability, we vary the number of objects from 100K to 1,000K. Other parameters
 381 are given their baseline values (Table 2). Figures 5, 6 and 7 show the effect of the number of objects.
 382 Three graphs are shown in a linear scale. For each algorithm, the runtime per iteration is
 383 approximately proportional to the number of objects.

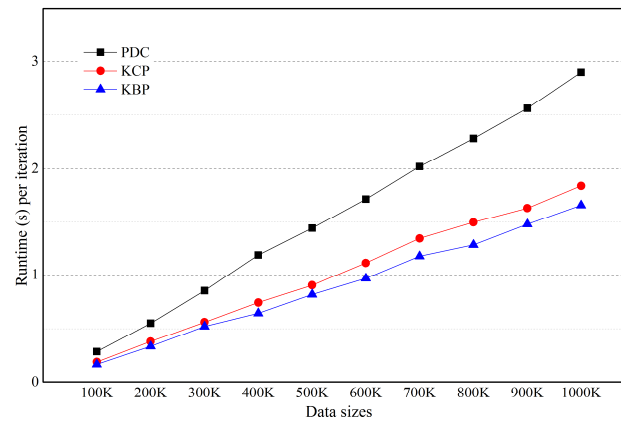
384 In Figure 5, KBP and KCP outperforms PDC. However, there is little difference in the
 385 performance between KBP and KCP. This is because if the categorical data is uniform distribution,
 386 most of the categorical data exist in the bitmap index. In this case, KBP is the same performance as
 387 KCP.



388
 389 **Figure 5.** Effect of the number of objects (numeric data and categorical data are on uniform
 390 distribution)

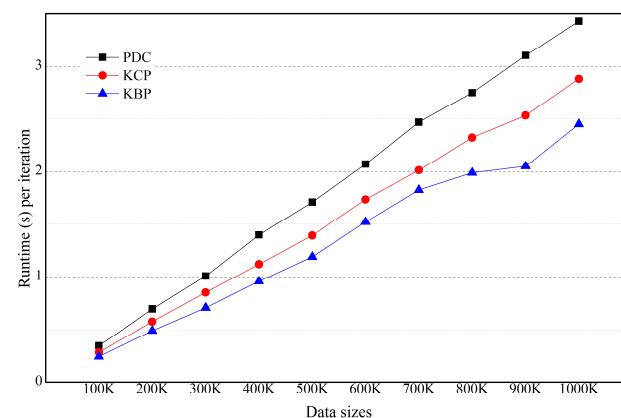
391
392
393
394
395
396
397
398

In Figure 6, KBP outperforms KCP and PDC. For example, KBP runs up to 1.1 and 1.75 times faster than KCP and PDC, respectively ($n=1,000K$). If categorical data is on a skewed distribution, KBP is effective for improving performance. As the size of the data increases, the difference in execution time increases. This is because as the data size increases, the amount of distance calculation increases, while at the same time the number of objects included in the cluster being pruned increases. In Figure 7, KCP outperforms PDC. Even if numeric data is on a Gaussian distribution, cell pruning is effective for improving performance.



399
400
401

Figure 6. Effect of the number of objects (numeric data is on uniform distribution and categorical data is on skewed distribution)



402
403
404

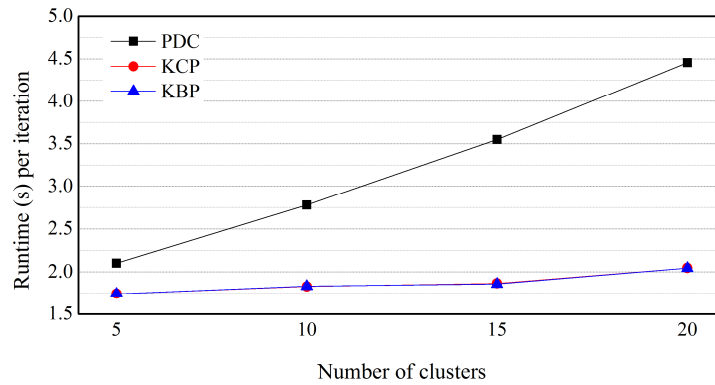
Figure 7. Effect of the number of objects (numeric data is on Gaussian distribution and categorical data is on skewed distribution)

405 5.3. Effects of the number of clusters

406
407
408
409
410
411
412
413

To confirm the effects of the number of clusters, we vary the number of clusters, k , from 5 and 20. Other parameters are kept at their baseline values (Table 2). Figures 8, 9 and 10 show the effects of the number of clusters. Three graphs are also shown in a linear scale. For each algorithm, the runtime is approximately proportional to the number of cluster.

In Figure 8, KBP and KCP also outperform PDC. However, there is also little difference in the performance between KBP and KCP. This is because if the categorical data is uniform distribution, most of the categorical data exist in the bitmap index like Figure 5.



414

415

416

Figure 8. Effect of the number of clusters (numeric data and categorical data are on uniform distribution)

417

418

419

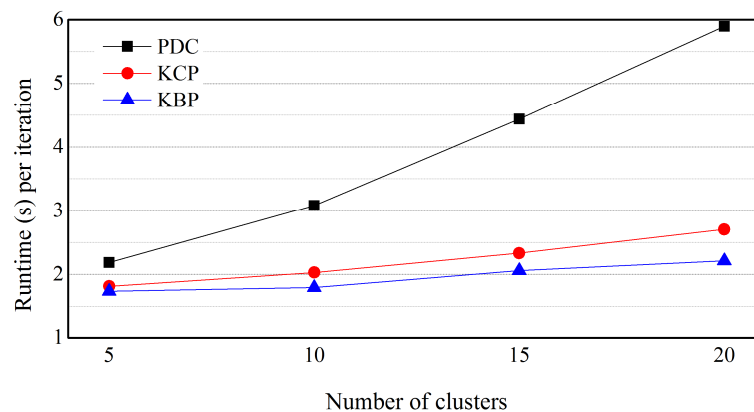
420

421

422

423

In Figure 9, KBP outperforms KCP and PDC. For example, KBP runs up to 1.13 and 1.71 times faster than KCP and PDC, respectively ($k = 10$). This result indicates that KBP is effective for improving performance even if categorical data is on a skewed distribution. As the number of cluster increases, the difference in execution time increases. This is also because as the data size increases, the amount of distance calculation increases, while at the same time the number of objects included in the cluster being pruned increases like Figure 6. In Figure 10, KCP also outperforms PDC. Even if numeric data is on a Gaussian distribution, KCP is also effective for improving performance.

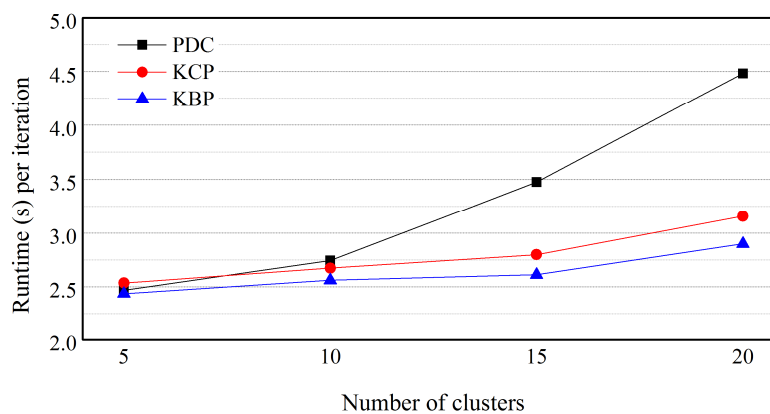


424

425

426

Figure 9. Effect of the number of clusters (numeric data and categorical data are on uniform distribution)



427

428

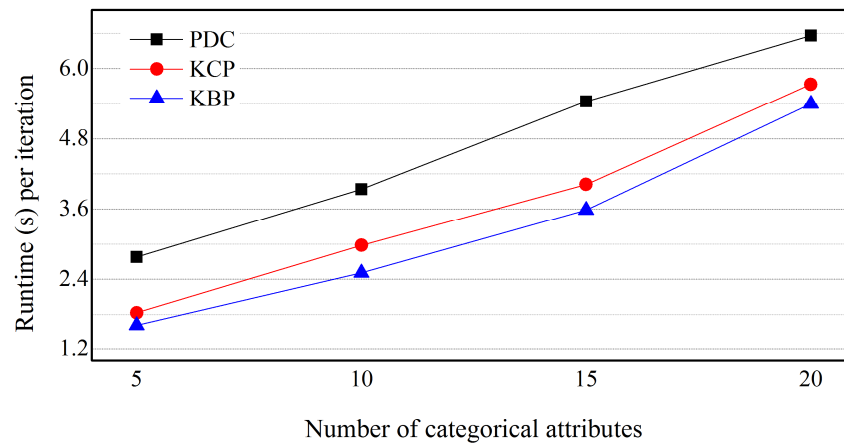
429

Figure 10. Effect of the number of clusters (numeric data is on Gaussian distribution and categorical data is on skewed distribution)

430

431 5.4. Effects of the number of categorical attributes

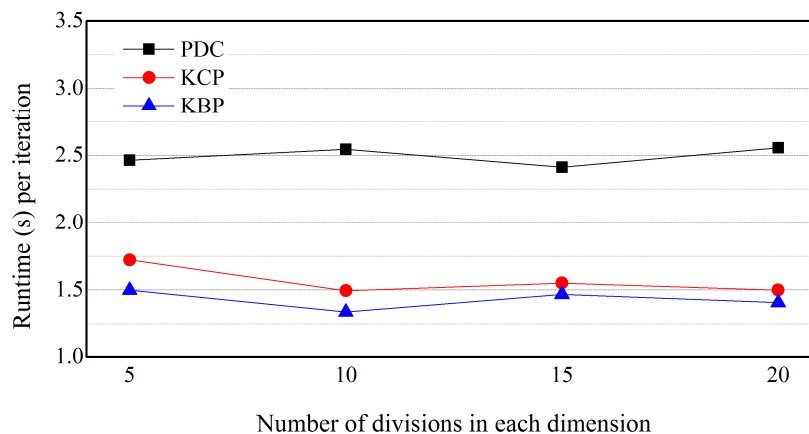
432 To confirm the effects of the number of categorical attributes, we vary the number of categorical
 433 attributes from 5 to 20. Other parameters are given their baseline values. Figure 11 shows the effect
 434 of the number of categorical attributes. The graph is also shown in a linear scale. For each algorithm,
 435 the runtime per iteration is approximately proportional to the number of categorical attributes. KBP
 436 outperforms KCP and PDC. For example, KBP runs up to 1.13 and 1.71 times faster than KCP and
 437 PDC, respectively ($m_c=5$). Even if the number of the categorical attributes increases, the difference
 438 between the execution time of KCP and PDC is kept almost constant. The reason is that KCP is based
 439 on numeric attributes and is not affected by the number of categorical attributes.
 440



441
 442 **Figure 11.** Effect of the number of categorical attributes on a skewed distribution

443 5.5. Effects of the size of cells

444 To confirm the effects of the size of cells, we vary the number of cells from 5 to 20. Other
 445 parameters are given their baseline values (Table 2). Figure 12 shows the effect of the number of
 446 divisions in each dimension. KBP and KCP outperform PDC. In Fig. 13, the horizontal axis is the
 447 number of divisions of each dimension. As the number of divisions increases, the size of the cell
 448 decreases. As the cell size gets smaller, the distance between the cell and cluster centers can be
 449 measured more finely. There is no significant difference in execution time according to the size of cell
 450 by each algorithm. This is because the distance calculation between the cell and the cluster centers is
 451 increased in proportion to the number of cells, and the bitmap index stored by the cell is also
 452 increased.



453
 454 **Figure 12.** Effect of the number of cells on a skewed distribution

455

456 **6. Conclusions**

457 In this paper we have propose an efficient grid-based k-prototypes algorithm, GK-prototypes,
 458 that improves performance for clustering spatial objects. We develop two pruning techniques, KCP
 459 and KBP. KCP which uses both maximum and minimum distance between cluster centers and a cell
 460 improves the performance than PDC. KBP is an extension of cell pruning for improving the efficiency
 461 of pruning in case that a categorical data is skewed. Our experimental results demonstrate that KCP
 462 and KBP outperforms PDC, and KBP outperforms KCP except for uniform distributions of categorical
 463 data. These results lead us to conclude that our grid-based k-prototypes algorithm can achieve better
 464 performance than the existing k-prototypes algorithm.

465 As data has grown exponentially and more complex recently, the traditional clustering
 466 algorithms have a great challenge to deal with these data. In future works, we may consider
 467 optimized pruning techniques of the k-prototypes algorithm in parallel processing environment.

468 **Author Contributions:** Conceptualization, Hong-Jun Jang and Byoungwook Kim; Methodology, Hong-Jun
 469 Jang; Software, Byoungwook Kim; Writing-Original Draft Preparation, Hong-Jun Jang and Byoungwook Kim;
 470 Writing-Review & Editing, Jongwan Kim; Supervision, Soon-Young Jung.

471 **Acknowledgments:** This research was supported by Basic Science Research Program through the National
 472 Research Foundation of Korea(NRF) funded by the Ministry of Education(No. NRF-2017R1D1A1B03034067) and
 473 by the National Research Foundation of Korea(NRF) grant funded by the Korea government(MSIT) (No. NRF-
 474 2016R1A2B1014013)

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

476 **References**

- 477 1. Sander, J.; Ester, M.; Kriegel, H.P.; Xu, X. Density-Based Clustering in Spatial Databases: The Algorithm
 478 GDBSCAN and Its Applications. *Data Mining and Knowledge Discovery* **1998**, *2*, 169–194.
- 479 2. Koperski, K.; Han, J.; Stefanovic, N. An Efficient Two-Step Method for Classification of Spatial Data. In
 480 Proceedings of the International Symposium on Spatial Data Handling (SDH'98), Vancouver, Canada, 1998,
 481 pp. 45–54.
- 482 3. Koperski, K.; Han, J. (1995). Discovery of Spatial Association Rules in Geographic Information Databases.
 483 In Proceedings of the 4th International Symposium on Advances in Spatial Databases (SSD'95), 1995, pp.
 484 47–66.
- 485 4. Ester, M.; Frommelt, A.; Kriegel, H.P.; Sander, J. Algorithms for Characterization and Trend Detection in
 486 Spatial Databases. In Proceedings of the Fourth International Conference on Knowledge Discovery and
 487 Data Mining (KDD'98), 1998, pp. 44–50.
- 488 5. Deren, L.; Shuliang, W.; Wenzhong, S.; Xinzhou, W. On Spatial Data Mining and Knowledge Discovery.
 489 *Geomatics and Information Science of Wuhan Univers* **2001**, *26*, 491–499.
- 490 6. Boldt, M.; Borg, A. A statistical method for detecting significant temporal hotspots using LISA statistics. In
 491 Proceedings of the Intelligence and Security Informatics Conference (EISIC), 2017 European, 2017.
- 492 7. Chainey, S.; Reid, S.; Stuart, N. When is a Hotspot a Hotspot? A Procedure for Creating Statistically Robust
 493 Hotspot Maps of Crime. In *Innovations in GIS 9: Socio-economic applications of geographic information*
 494 *science*, Kidner, D; Higgs, G; White, S, Eds.; Taylor & Francis: London, UK, 2002; pp. 21–36.
- 495 8. Murray, A.; McGuffog, I.; Western, J.; Mullins, P. Exploratory spatial data analysis techniques for
 496 examining urban crime. *The British Journal of Criminology* **2001**, *41*, pp. 309–329.
- 497 9. Chainey, S.; Tompson, L.; Uhlig, S. The Utility of Hotspot Mapping for Predicting Spatial Patterns of Crime.
 498 *Security Journal* **2008**, *21*, pp. 4–28.
- 499 10. Di Martino, F.; Sessa, S. The extended fuzzy C-means algorithm for hotspots in spatio-temporal GIS. *Expert*
 500 *Systems with Applications* **2011**, *38*, pp. 11829–11836.
- 501 11. Di Martino, F.; Sessa, S.; Barillari, U.E.S.; Barillari, M.R. Spatio-temporal hotspots and application on a
 502 disease analysis case via GIS. *Soft Computing* **2014**, *18*, pp. 2377–2384.
- 503 12. Mullner, R.M.; Chung, K.; Croke, K. G.; Mensah, E. K. Geographic information systems in public health
 504 and medicine. *Journal of Medical Systems* **2004**, *28*, pp. 215–221.
- 505 13. Polat, K. Application of attribute weighting method based on clustering centers to discrimination of
 506 linearly non-separable medical datasets. *Journal of Medical Systems* **2012**, *36*, pp. 2657–2673.

- 507 14. Wei, C.K.; Su, S.; Yang, M.C. Application of data mining on the development of a disease distribution map
508 of screened community residents of taipei county in Taiwan. *Journal of Medical Systems* **2012**, *36*, pp. 2021–
509 2027.
- 510 15. Huang, Z. Clustering large data sets with mixed numeric and categorical values. In Proceedings of the First
511 Pacific Asia Knowledge Discovery and Data Mining Conference, 1997, pp. 21–34.
- 512 16. Kim, B. A Fast K-prototypes Algorithm Using Partial Distance Computation. *Symmetry* **2017**, *9*,
513 doi:10.3390/sym9040058
- 514 17. Goodchild, M. Geographical information science. *International Journal of Geographic Information Systems*
515 **1992**, *6*, pp. 31–45.
- 516 18. Xiaoyun, C.; Yi, C.; Xiaoli, Q.; Min, Y.; Yanshan, H. PGMCLU: a novel parallel grid-based clustering
517 algorithm for multi-density datasets. In: 1st IEEE symposium on web society, 2009 (SWS'09), Lanzhou,
518 2009, pp 166–171.
- 519 19. Wang, W.; Yang, J.; Muntz, R. R. STING: A Statistical Information Grid Approach to Spatial Data Mining.
520 In the 23rd International Conference on Very Large Data Bases (VLDB'97), 1997, pp. 186–195.
- 521 20. Agrawal, R.; Gehrke, J.; Gunopulos, D.; Raghavan, P. Automatic Subspace Clustering of High Dimensional
522 Data for Data Mining Applications. In Proceedings of the ACM SIGMOD International Conference on
523 Management of Data, ACM Press, 1998, pp. 94–105.
- 524 21. Chen, X., Su, Y., Chen, Y., & Liu, G. GK-means: An Efficient K-means Clustering Algorithm Based On Grid.
525 In Computer Network and Multimedia Technology (CNMT 2009) International Symposium, 2009.
- 526 22. Ji, J.; Pang, W.; Zheng, Y.; Wang, Z.; Ma, Z.; Zhang, L. A Novel Cluster Center Initialization Metho
527 d for the k-Prototypes Algorithms using Centrality and Distance. *Applied Mathematics & Information*
528 *Sciences* **2015**, *9*, pp. 2933–2942.
- 529 23. Zavadskas, E.K.; Antucheviciene, J.; Vilutiene, T.; Adeli, H. Sustainable Decision Making in Civil Enginee
530 ring, Construction and Building Technology. *Sustainability* **2018**, *10*, 14.
- 531 24. Hersh, M.A. Sustainable Decision Making: The Role of Decision Support systems. *IEEE Transactions on Sy*
532 *stems, Man, and Cybernetics-Part C: Applications and Reviews* **1999**, *29*, 3, pp. 395-408.
- 533 25. Morik, K.; Bhaduri, K.; Kargupta, H. Introduction to data mining for sustainability. *Data Mining and*
534 *Knowledge Discovery* **2012**, *24*, 2, pp. 311-324.
- 535 26. Aissi, S.; Gouider, M.S.; Sboui, T.; Said, L.B. A spatial data warehouse recommendation approach:co
536 nceptual framework and experimental evaluation. *Human-centric Computing and Information Sciences* **2**
537 **015**, *5*, 30.
- 538 27. Kim, J.-J. Spatio-temporal Sensor Data Processing Techniques. *Journal of Information Processing System*
539 *s* **2017**, *13*, 5, pp. 1259-1276.