

# Evaluation of different clustering techniques in classifying the vegetable growing panchayats of Ernakulam district, Kerala

## ABSTRACT

The goal of this study was to evaluate different clustering techniques in classifying the vegetable growing locations of Ernakulam (EKM) district of Kerala so that same nutrient recommendation can be prescribed for panchayats coming under the same cluster. Hierarchical clustering (HC) and  $K$ -means clustering were performed to group the panchayats based on soil fertility status and thereafter comparison of various clustering procedures was done using Davies – Bouldin (DB) index. Different dissimilarity measures- Euclidean, squared Euclidean, Chebychev distance and Mahalanobis  $D^2$  were determined and single linkage, complete linkage and average linkage methods were adopted under these measures. The results revealed that Mahalanobis  $D^2$  was the better clustering procedure with seven clusters (DB index: 0.120) followed by average linkage method under Euclidean distance (DB index: 0.306) with seven clusters. Manjapra and Keerampara panchayats remained as individual clusters. Keerampara had strongly acidic soils (pH -5.17) with high available Mg (158 mg kg<sup>-1</sup>) while Manjapra soils had low Mg availability (19 mg kg<sup>-1</sup>) and high S content (57 mg kg<sup>-1</sup>). Kakkad, Kalady and Vengoor came under cluster I which possessed approximately same EC (0.15-0.19 dS m<sup>-1</sup>), OC (2-2.4%) and Mg (71-73 mg kg<sup>-1</sup>) content. Chengamanadu and Vengola came under cluster III while Ayyampuzha and Mudakkuzha belong to cluster IV.

*Keywords: Hierarchical clustering, Davies – Bouldin index, Mahalanobis  $D^2$ , Euclidean distance, Average linkage.+*

## 1. INTRODUCTION

Kerala, the God's own country is blessed with biological diversity and soil heterogeneity. Ernakulam is a district situated in the central region of Kerala and rice is the major crop cultivated in the wet lands of the district. In addition to rice, vegetables, pineapple and nutmeg are also cultivated in EKM. Plant growth is highly dependent on the fertility status of soil prevailing in the region [1]. Plants need primary nutrients (Nitrogen (N), Phosphorus (P), Potassium (K)), secondary nutrients (Calcium (Ca), Magnesium (Mg), Sulphur (S)) and micro nutrients (Boron (B), Iron (Fe), Manganese (Mn), Zinc (Zn), Copper (Cu)) in adequate quantities for growth and development and they are absorbed from the soil. Other soil parameters viz. pH and Electrical conductivity (EC) also affect plant growth and nutrient availability [2].

Systematic testing of soil nutrients including the micro nutrients and prescribing recommendation is essential for soil health management. The soil test data available for each panchayat can be used for classifying them into various groups based on similarity in soil parameters. Cluster analysis, one of the multivariate methods is suitable in classifying objects based on similarity measures [3]. Clustering of panchayats based on their soil fertility status can be accomplished by means of cluster analysis [4]. Panchayats coming under the same cluster have similar soil characteristics and that under different clusters have dissimilar soil properties. Different clustering procedures are available which vary according to the distance measures selected for clustering. Same recommendations can be given to those panchayats which come under the same cluster.

Cluster analysis is commonly applied in the field of plant breeding and genetics where varieties/genotypes can be grouped based on their quantitative characters [5]. According to [6], results of cluster analysis changes when different methods of clustering are used. Cluster analysis

45 can be used to partition sites based on soil characteristics [7]. The land use effect on soil chemical  
46 and microbial properties can be determined with the use of cluster analysis [8]. Application of cluster  
47 analysis on soil geochemical data helped to determine the spatial distribution of elements [9].  
48 Evaluation of different clustering techniques is possible with a measure; Davies- Bouldin index which  
49 identifies the cluster that best fit the data [10].

50 This study is an attempt to classify the panchayats of Ernakulam district based on the soil  
51 fertility status. Among the various clustering procedures, clusters that best fit the data are also  
52 identified.

## 53 **2. MATERIALS AND METHODS**

54 The present study is based on the data on thirteen soil fertility parameters of different vegetable  
55 growing locations (panchayats) of Ernakulam district of Kerala. Soil samples collected from different  
56 panchayats of Ernakulam district during the year 2016-2017 analysed by Department of Soil science  
57 & Agricultural Chemistry, College of Agriculture, Vellayani, Kerala and the data maintained was  
58 utilized for the present study. Sample collection was done by the farmers themselves from their own  
59 vegetable growing plots at a depth of 15 cm using spade. From each sampling plot, about 10-15  
60 random samples were collected, mixed and reduced to 0.5 Kg by the method of quartering. It was  
61 observed that panchayats show variation in soil properties with respect to cropping patterns and  
62 cultivation practices [11].

63 The data on thirteen soil fertility parameters of 17 panchayats of Ernakulam viz., electro chemical  
64 parameters (pH and Electrical conductivity (EC)), Oxidisable Organic Carbon (OC), Phosphorus (P),  
65 Potassium (K), secondary nutrients (Calcium (Ca), Magnesium (Mg), Sulphur (S)) and micro nutrients  
66 (Boron (B), Iron (Fe), Manganese (Mn), Zinc (Zn), Copper (Cu)) were available. Each panchayat is  
67 having different sample sizes and altogether sample size comes around 583.

### 68 **2.1 Cluster analysis (CA)**

69 Cluster analysis is a multivariate technique used to group individuals or objects based on their  
70 several characteristics [12]. There should be homogeneity within groups or clusters and heterogeneity  
71 between groups. First we measure the distance between objects based on their multiple characters  
72 which is otherwise called as similarity or dissimilarity measures. Based on the similarity/dissimilarity  
73 measures, clusters are formed later.

#### 74 **2.1.1 Distance measures**

75 Euclidean distance, squared Euclidean distance, Chebychev distance and Mahalanobis  $D^2$   
76 are the distance measures used for the present work.

##### 77 **2.1.1.1 Euclidean distance**

78 It is the geometrical distance between two objects in the multidimensional space. It is  
79 calculated as,

$$80 \quad d(x, y) = \sqrt{(X - Y)'(X - Y)}$$

81 where  $X$  and  $Y$  are the ' $p$ ' dimensional vector of observations and  $X = (X_1, X_2, \dots, X_p)$  and  $Y =$   
82  $(Y_1, Y_2, \dots, Y_p)$ . Euclidean distance is one of the most commonly used distance measures.

##### 83 **2.1.1.2 Squared Euclidean distance**

84 It is the square of Euclidean distance and is used to put weights for those objects which are  
85 farther apart.

86 
$$E_{ij} = \sum_{k=1}^p (X_{ik} - X_{jk})^2$$

87 **2.1.1.3 Mahalanobis  $D^2$  statistics**

88 A measure for group distance based on multiple characters was given by Mahalanobis, 1936. With  $x_1,$   
 89  $x_2, x_3, \dots, x_p$  as multiple measurements available on each individual and  $d_1, d_2, d_3, \dots, d_p$  as  $\bar{x}_1^1 - \bar{x}_1^2,$   
 90  $\bar{x}_2^1 - \bar{x}_2^2, \dots, \bar{x}_p^1 - \bar{x}_p^2,$  respectively, being the difference in the means of two populations,  
 91 Mahalanobis  $D^2$  statistics is defined as follows:

92 
$$D^2 = b_1 d_1 + b_2 d_2 + \dots + b_p d_p = \sum_i \sum_j (\bar{X}_1 - \bar{X}_2)' W^{-1} (\bar{X}_1 - \bar{X}_2)$$

93 Where  $\bar{x}_1^1$  is the mean value of 1<sup>st</sup> character in the first population and  $\bar{x}_1^2$  is the mean of the 1<sup>st</sup>  
 94 character in the second population. Here, the  $b_i$  values are to be estimated such that the ratio of  
 95 variance between the populations to the variance within the population is maximized. In terms of  
 96 variances and covariances, the  $D^2$  distance between object 1 and object 2 can be obtained as;

97 
$$D^2 = (\bar{X}_1 - \bar{X}_2)' W^{-1} (\bar{X}_1 - \bar{X}_2)$$

98 Where  $W^{-1}$  is the inverse of variance covariance matrix,  $\bar{X}_1$  is the mean of first population,  $\bar{X}_2$  is  
 99 the mean of second population. It is used for quantitative data.

100 **2.1.1.4 Chebychev's distance**

101 This is a dissimilarity measure is based on the assumption that two objects are different if  
 102 they differ in any one of the characteristics and is calculated as,

103 
$$C_{ij} = \text{Max } |X_{ik} - X_{jk}|$$

104 **2.1.2 Clustering techniques**

105 These are the procedures used for clubbing together of similar objects into different clusters. There  
 106 are different techniques for clustering *i.e.* Hierarchical technique and *K*-means clustering [13].

107 **2.1.2.1 Hierarchical clustering**

108 Hierarchical clustering assumes each of the '*n*' objects as individual clusters initially. Similar  
 109 objects are combined together in successive fusions or dissimilar objects are divided in successive  
 110 divisions in this clustering technique. There are different methods for linking the objects in different  
 111 clusters. Single linkage (nearest neighbour) method, complete linkage (farthest neighbour) method  
 112 and Tocher's method are common.

113 In single linkage, two individuals having minimum distance forms the first cluster. In next step, a third  
 114 individual is joined with the initial cluster or another two nearer individuals are clustered together to  
 115 form the second cluster. This is determined if the distance from the third individual to the first cluster is  
 116 shorter than the distance between the two nearer individuals. Two objects having maximum distance  
 117 between them constitute two groups. Next object either join one of the previous clusters or form its  
 118 own cluster.

119 Tocher's method of clustering objects makes use of Mahalanobis  $D^2$  statistic.  $D^2$  values are  
 120 arranged in ascending order and the two individuals having smallest distance between them is  
 121 selected as the first cluster. Tocher suggested a cut off value which is equal to maximum among the  
 122 minimum  $D^2$  values. Addition of a third object is determined in such a way that average  $D^2$  distance is  
 123 less than the cut off value. Clustering continues until all the objects are included in one of the clusters  
 124 [14].

### 125 **2.1.2.2 K -means clustering**

126 K-means clustering is a technique where we define some pre-specified number of clusters for  
127 grouping 'n' objects. Only condition is that clusters are formed with minimum variance within clusters  
128 and maximum variance between clusters.

### 129 **2.1.3 Cluster validity index**

130 Davies-Bouldin index is one of the cluster validity measures to evaluate the clusters based on  
131 their compactness and separation from each other. Let  $X_1, X_2, \dots, X_c$  are the clusters and  $\Delta x_i, \delta(x_i, x_j)$   
132 represents the intra cluster distance of  $i^{\text{th}}$  cluster and inter cluster distance between  $i^{\text{th}}$  and  $j^{\text{th}}$  cluster  
133 respectively. Then DB index is defined as,

$$134 \quad DB = \frac{1}{c} \sum_{i=1}^c \max_{i \neq j} \frac{(\Delta x_i + \Delta x_j)}{\delta(x_i, x_j)}$$

135

136 DB index values should be less for good clustering *i.e.* small index values corresponds to clusters  
137 which are compact and well separated. Clustering algorithms that minimize the DB index values give  
138 optimum number of clusters.

## 139 **3. RESULTS AND DISCUSSION**

140 In order to classify the various panchayats based on the soil fertility status in Ernakulam  
141 district, cluster analysis was adopted [15]. Different clustering algorithms such as K-means clustering  
142 and hierarchical agglomerative clustering were adopted in the present study. Various dissimilarity  
143 measures such as Euclidean distance, squared Euclidean, Mahalanobis  $D^2$  and Chebychev distance  
144 were estimated for clustering the panchayats. Various clustering procedures like single linkage,  
145 complete linkage and average linkage were followed under hierarchical clustering to generate clusters  
146 consisting of panchyats having similar soil fertility properties [16]. Tocher's clustering procedure was  
147 also employed to generate clusters using Mahalanobis  $D^2$  distance measure. Comparison of various  
148 measures and clustering algorithms were completed using SPSS package. The study conducted by  
149 [17] was in accordance with the present study as different clustering techniques provided different  
150 number of clusters.

### 151 **3.1 Hierarchical clustering of panchayats**

#### 152 **3.1.1 Squared Euclidean distance**

153 Squared Euclidean distance was selected for clustering the panchayats and the dissimilarity  
154 matrix was determined. Single linkage, complete linkage procedures were practiced along with  
155 average linkage method for squared Euclidean distance and the clustering of panchayats in each  
156 case is given in Table 1.

157 There were seven clusters under single linkage and average linkage clustering procedure and  
158 six clusters for complete linkage. Single linkage was another hierarchical clustering procedure and it  
159 did not provide proper clusters of panchayats in EKM district.

160 It is evident from the Table 1 that Ayyampuzha, Chengamanadu, Keerampara, Thirumaradi  
161 and Thuravur formed as single clusters when single linkage was adopted. Five panchayats out of 17  
162 were in individual clusters. Only two clusters were having more than one panchayat in it and all other  
163 panchayats remained as individual clusters.

164 The panchayats Ayyampuzha, Chengamanadu, and Thuravur remained as single clusters  
 165 under complete linkage and average linkage methods and Kadungallor, Kalady and Puthenvelikkra  
 166 were also in same cluster under these clustering procedures. There were six clusters under complete  
 167 linkage and seven clusters under average linkage. Ayyampuzha, Chengamanadu and Thuravur stood  
 168 as individual clusters,(Table 1).

169 **Table 1. Clustering of panchayats in EKM based on single linkage, complete linkage and**  
 170 **average linkage method (Squared Euclidean distance).**

Cluster	Panchayats coming under each cluster		
	Single linkage	Complete linkage	Average linkage
I	Ayyampuzha	Ayyampuzha	Ayyampuzha
II	Kadungalloor, Kalady, Puthenvelikkara	Kadungalloor, Kalady, Puthenvelikkara	Kadungalloor, Kalady, Puthenvelikkara
III	Thuravur	Thuravur	Thuravur
IV	Chengamanadu	Chengamanadu	Chengamanadu
V	Keerampara,	Keerampara, Thirumaradi, Mudakkuzha, Nedumbassery, Piravom	Keerampara, Mudakkuzha, Thirumaradi
VI	Mudakkuzha, Kakkad, Vengola, Manjapra, Vengoor, Nedumbassery, Piravom, Pampakuda, Pothanikkad	Kakkad, Vengola, Manjapra, Vengoor, Pampakuda, Pothanikkad	Kakkad, Vengola, Manjapra, Vengoor
VII	Thirumaradi		Nedumbassery, Piravom, Pampakuda, Pothanikkad

171

172 **3.1.2 Euclidean distance**

173 Euclidean distance was also selected as the distance measure and single linkage, complete  
 174 linkage and average linkage procedures were performed under this distance measure. Average  
 175 linkage using Euclidean distance as the distance measure gave the same clustering pattern as that of  
 176 squared Euclidean (Table 2). Unlike squared Euclidean, average linkage and complete linkage using  
 177 Euclidean distance resulted in seven clusters and single linkage method had six clusters.  
 178 Ayyampuzha, Chengamanadu, Keerampara and Thuravur were formed as clusters with single  
 179 panchayat in it and hence intra cluster distance was zero. Except Keerampara which was a single  
 180 cluster for single linkage method only, Ayyampuzha, Chengamanadu and Thuravur were also found  
 181 to be single clusters when complete and average linkage methods were used. Similar findings were  
 182 put forwarded by [18].

Cluster	Panchayats coming under each cluster
---------	--------------------------------------

	Sinkle linkage	Complete linkage	Average linkage
I	Ayyampuzha	Ayyampuzha	Ayyampuzha
II	Thuravur	Thuravur	Thuravur
III	Chengamanadu	Chengamanadu	Chengamanadu
IV	Kadungalloor, Kalady, Puthenvelikkara	Kadungalloor, Kalady, Puthenvelikkara	Kadungalloor, Kalady, Puthenvelikkara
V	Keerampara	Keerampara, Mudakkuzha, Thirumaradi, Nedumbassery, Piravom	Keerampara, Mudakkuzha, Thirumaradi
VI	Nedumbassery, Piravom, Kakkad, Mudakkuzha, Pampakuda, Pothanikkad, Vengola, Vengoor, Manjapra, Thirumaradi	Kakkad, Vengola, Manjapra, Vengoor, Pampakuda, Pothanikkad	Kakkad, Vengola, Manjapra, Vengoor
VII		Kadungalloor, Kalady, Puthenvelikkara	Nedumbassery, Piravom, Pampakuda, Pothanikkad,

183 **Table 2. Clustering of panchayats in EKM based on single linkage, complete linkage and**  
184 **average linkage method (Euclidean distance).**

185

186

187 Kadungalloor, Kalady and Puthenvelikkara came under a cluster when all these methods  
188 were practiced. This revealed that these three panchayats had almost similar soil characteristics.  
189 Similarly, Kakkad, Vengola, Manjapra and Vengoor were the panchayats which formed into a cluster  
190 and it was inferred that these panchayats also had similar soil properties.

191 Clusters obtained through Euclidean distance were similar as that of squared Euclidean distance.  
192 Ayyampuzha, Chengamanadu and Thuravur retained as individual clusters with intra cluster distance  
193 zero.

### 194 3.1.3 Mahalanobis $D^2$ distance

195 Panchayats coming under each cluster as per Mahalanobis  $D^2$  distance is given in Table 3. This  
196 distance measure provided seven clusters when Tocher's clustering method was adopted. Manjapra  
197 and Keerampara were the panchayats which retained as single cluster with intra cluster distance zero  
198 when Mahalanobis  $D^2$  was selected as the distance measure. This result is in line with the results of  
199 the research work done by [5] on clustering of rice genotypes.

200

201

202

203 **Table 3. Clustering of panchayats in EKM based on Tocher's method**

Cluster	Panchayaths coming under each cluster	204
I	Kakkad, Kalady, Vengoor	205
II	Kadungalloor, Pothanikkad, Puthenvelikkara, Pampakuda, Thirumaradi	206
III	Chengamanadu, Vengola	207
IV	Ayyampuzha, Mudakkuzha	208
V	Nedumbassery, Piravom	209
VI	Manjapra	209
VII	Keerampara	210

210

### 211 3.1.4 Chebychev distance

212 Clustering of panchayats based on average linkage, single linkage and complete linkage  
 213 methods using Chebychev distance as the dissimilarity measure is given below in Table 4. There  
 214 were 12 clusters into which panchayats were grouped when single linkage method was adopted.  
 215 Other than Ayyampuzha, Chengamanadu and Thuravur which were separate clusters in all the  
 216 clustering procedures, Manjapra, Vengoor, Thirumaradi, Puthenvelikkara, Mudakkuzha and  
 217 Keerampara stood as individual clusters. Kadungalloor, Kalady and Puthenvelikkara were formed as  
 218 a single cluster inferring that these three panchayats had almost similar soil nutrient status. Complete  
 219 linkage and average linkage gave seven clusters. Mudakkuzha, Keerampara and Thirumaradi came  
 220 under the same cluster when both average and complete linkage methods were adopted.

221 **Table 4. Clustering of panchayats in EKM based on Chebychev distance under different**  
 222 **methods**

Cluster	Panchayats coming under each cluster			223
	Single linkage	Complete linkage	Average linkage	
I	Ayyampuzha	Ayyampuzha	Ayyampuzha	224
II	Thuravur	Thuravur	Thuravur	225
III	Kadungalloor, Kalady	Kadungalloor, Kalady, Puthenvelikkara	Kadungalloor, Kalady, Puthenvelikkara	226
IV	Chengamanadu	Chengamanadu, Nedumbassery, Piravom	Chengamanadu	227
V	Mudakkuzha	Mudakkuzha, Keerampara, Thirumaradi	Mudakkuzha, Keerampara, Thirumaradi	228 229
VI	Vengola, Kakkad,	Vengola, Kakkad, Manjapra, Vengoor	Vengola, Kakkad, Manjapra, Vengoor	230
VII	Pothanikkad, Pampakuda, Nedumbassery, Piravom	Pothanikkad, Pampakuda	Pothanikkad, Pampakuda Nedumbassery, Piravom	231 232
VIII	Vengoor			233
IX	Puthenvelikkara			234
X	Keerampara			234
XI	Manjapra			235
XII	Thirumaradi			235

236 Chebychev distance also provided somewhat similar clustering patterns as that of squared Euclidean  
 237 and Euclidean distance. Only difference was in the number of clusters formed when single linkage

238 method used under the three distance measures. Chebychev distance gave 12 clusters while  
 239 Euclidean and squared Euclidean gave 6 and 7 clusters respectively.

### 240 3.2 K-means clustering

241 K-means clustering defines the number of clusters in advance. In the present study, number  
 242 of clusters was taken as 5, 6 and 7 in K-means clustering. When  $k = 5$ ,  $K = 6$  and  $K = 7$  panchayats of  
 243 Ernakulam was classified into five, six and seven groups respectively and are presented in Table 5.

244 Only Ayyampuzha was as an individual cluster for  $K = 5, 6$  and  $7$ . Thuravur became a single cluster  
 245 when  $k = 6$ . Nedumbassery, Piravom and Chengamanadu came under cluster II for  $K = 5, 6$  and  
 246  $7$ . Kadungalloor, Kalady and Puthovelikkara belonged to same cluster which suggested that they had  
 247 similar soil properties. The same clustering procedure was adopted by [7] for partitioning locations  
 248 based on soil properties.

249 Thuravur, which was in the same cluster with other panchayats, now formed as a new individual  
 250 cluster when  $K = 6$  was selected. Panchayats were again redistributed among different clusters when  
 251 K-means clustering with  $K = 7$  were performed. Now, Ayyampuzha, Keerampara and Thuravur were  
 252 in different clusters with intra cluster distance zero i.e. those cluster had single panchayats as the  
 253 entity.

254 **Table 5. Clustering of panchayats in EKM based on K-means clustering**

Cluster	Panchayats coming under each cluster			255
	K =5	K =6	K =7	256
I	Ayyampuzha	Ayyampuzha	Ayyampuzha	257
II	Chengamanadu, Nedumbassery, Piravom	Chengamanadu, Nedumbassery, Piravom	Chengamanadu, Nedumbassery, Piravom	258 259
III	Kadungalloor, Kalady, Puthenvelikkara	Kadungalloor, Kalady, Puthenvelikkara	Kadungalloor, Kalady, Puthenvelikkara	260 261 262
IV	Kakkad, Manjapra, Pampakuda, Pothanikkad, Thuravur, Vengola, Vengoor	Kakkad, Manjapra, Pampakuda, Pothanikkad, Vengola, Vengoor	Kakkad, Manjapra, Pothanikkad, Vengola, Vengoor	263 264 265
V	Keerampara, Mudakkuzha, Thirumaradi	Keerampara, Mudakkuzha, Thirumaradi	Mudakkuzha, Pampakuda, Thirumaradi	266 267
VI		Thuravur	Thuravur	268
VII			Keerampara	269

271 It was concluded that Ayyampuzha, Thuravur and Chengamanadu were the panchayats with different  
 272 soil characteristics and did not form any group with other panchayats based on the soil fertility



273 parameters. In most of the clustering procedures adopted, Kadungalloor, Kalady and Puthenvelikkara  
 274 came under a cluster based on their soil properties which indicated that these three panchayats had  
 275 almost similar soil properties.

### 276 3.3. Cluster validity index- panchayats of Ernakulam

277 Being unsupervised procedure, cluster analysis need evaluation of the results of different clustering  
 278 procedures. Cluster validity means identifying the clusters that best fit to the given data. Davies-  
 279 Bouldin index was one of such measures used for cluster validation which was used in the present  
 280 study (Table 6). Comparatively low values are preferred for good clustering procedure in this method.

281 **Table 6. Cluster validity index for different clustering procedures (EKM)**

Distance measure	Linkage method	DB index score	282
Squared Euclidean	Average linkage	0.412	283
	Complete linkage	0.427	284
Euclidean distance	Average linkage	<b>0.306</b>	285
	Complete linkage	0.894	286
Chebychev distance	Average linkage	0.458	287
	Complete linkage	0.383	288
Mahalanobis $D^2$		<b>0.120</b>	288
$K$ – means clustering	$K=5$	0.566	289
	$K=6$	0.467	290
	$K=7$	0.497	291

292

293 Single linkage method was not considered as most of the panchayat was retained as single  
 294 clusters with intra cluster distance zero. A comparison between average linkage and complete linkage  
 295 was carried out using the D-B index [19]. D- B index should be less for the optimum clustering pattern  
 296 and here Mahalanobis  $D^2$  was found to be the best clustering measure followed by average linkage  
 297 method with Euclidean as the distance measure. Cluster means were determined for Mahalanobis  $D^2$   
 298 and the results are given in Table 7.

299

300

301

302

303

304

305

306 **Table 7. Cluster mean based on Mahalanobi's  $D^2$  (EKM)**

Cluster no.	pH	EC	OC	P	K	Ca	Mg	S	B	Fe	Mn	Zn	Cu
I	5.2 97	0.1 44	1.9 59	47.0 67	298.5 62	530.0 51	71.95 3	22.3 64	0.5 46	85.71 3	26.0 12	3.6 73	2.8 05
II	4.8 80	0.1 29	1.6 68	49.5 22	191.9 96	520.9 96	29.11 4	28.5 79	0.7 43	60.60 5	17.9 56	2.1 22	4.3 46
III	4.6 40	0.1 97	1.5 73	86.1 73	231.7 98	414.3 39	55.54 1	24.9 20	0.9 06	35.64 9	12.2 52	1.3 49	4.2 66
IV	5.2 07	0.1 51	1.9 11	56.7 47	346.9 30	522.5 71	108.9 95	13.5 25	0.7 17	76.64 9	29.6 85	2.4 36	3.3 13
V	5.0 83	0.1 03	1.3 94	55.6 13	165.1 48	369.0 00	80.50 0	20.9 48	0.8 63	38.03 5	13.6 10	1.7 97	1.5 05
VI	5.2 84	0.2 84	1.7 89	82.7 74	209.0 16	490.5 26	19.68 4	57.6 53	1.0 76	115.7 90	27.5 00	1.6 21	2.0 11
VII	5.1 71	0.0 66	1.9 28	31.2 03	139.1 12	274.6 55	41.17 2	21.2 47	1.0 15	158.2 41	21.7 64	5.7 09	5.2 22

307

308 Keerampara and Manjapra remained as separate panchayats when Mahalanobis  $D^2$  was  
309 practiced. Even though both panchayats come under the agro ecological unit 3.1 (southern and  
310 central foot hills), they were not found together under a cluster as they would be having dissimilar soil  
311 characteristics [20]. Keerampara soils were strongly acidic with comparatively lower available Ca  
312 content (274 mg kg<sup>-1</sup>) and high Fe status (158 mg kg<sup>-1</sup>). Some parts of Keerampara had proximity  
313 with water bodies. Nedumbassery and Piravom belonged to cluster V with EC (0.10 dS m<sup>-1</sup>), OC  
314 (1.39 %), K (165.14 kg ha<sup>-1</sup>) and S (20.94 mg kg<sup>-1</sup>). Cluster IV comprised Ayyampuzha and  
315 Mudakkuzha panchayats with an average pH (5.2), EC (0.15), OC (1.91 %) and K (346.93 kg ha<sup>-1</sup>).  
316 Kadungalloor, Pothanikkad, Puthenvelikkara, Pampakuda and Thirumaradi clubbed together to form a  
317 cluster with 1.6 per cent OC, P (49.5 kg ha<sup>-1</sup>) and K (191.99 kg ha<sup>-1</sup>). Kakkad, Kalady and Vengoor  
318 came under cluster I as they had approximately the same EC (0.15-0.19 dS m<sup>-1</sup>), OC (2-2.4%) and  
319 Mg (71-73 mg kg<sup>-1</sup>). Chengamanadu which was reported as an individual cluster in all the clustering  
320 methods adopted, came along with Vengola under Tocher's method. OC ranged from 1.55 to 1.58  
321 per cent, K (216.88-246.71 kg ha<sup>-1</sup>), Mg (54-56 mg kg<sup>-1</sup>) and B (0.89-0.91 mg kg<sup>-1</sup>) in these  
322 panchayats.

323 For all the clustering methods used except Mahalanobis  $D^2$ , Ayyampuzha retained as a single  
324 cluster. Hierarchical clustering also provided information that Ayyampuzha did not form any group  
325 with other panchayats. Ayyampuzha had moderately acidic soils (pH - 5.8) with high available K (570  
326 kg ha<sup>-1</sup>) while Chengamanadu had high available P (113 kg ha<sup>-1</sup>) and low Ca content (300 mg kg<sup>-1</sup>).  
327 Ayyampuzha comes under the agro ecological unit 4.1 (Southern high hills) and lies near to  
328 waterbodies. Thuravur was deficient in available B (0.28mg kg<sup>-1</sup>). Kadungalloor, Kalady and  
329 Puthenvelikkara came under cluster IV and the soils in this cluster recorded with high K (342.82 kg  
330 ha<sup>-1</sup>) and Ca (619.82 mg kg<sup>-1</sup>). The cluster V comprised of Keerampara, Mudakkuzha and  
331 Thirumarady and the soils were strongly acidic with comparatively lower P (33.84 kg ha<sup>-1</sup>) status. The  
332 soils of panchayats in cluster VI were moderately acidic with sufficient quantity of soil nutrients.  
333 Cluster means calculated for different clusters based on Euclidean distance with average linkage  
334 method is presented in Table 8.

335

336

337

338

**Table 8. Cluster mean based on Euclidean distance (EKM)**

Cluster no.	pH	EC	OC	P	K	Ca	Mg	S	B	Fe	Mn	Zn	Cu
I	5.8 4	0.1 5	1.7 3	63.72	570.2 1	670.7 0	109.7 7	8.12	0.5 0	19.99	16.4 0	2.3 0	1.1 7
II	5.0 9	0.1 5	1.7 4	41.61	112.9 8	616.4 4	33.68	16.0 5	0.2 8	37.77	29.1 2	6.5 1	4.3 4
III	4.2 4	0.2 5	1.5 6	113.5 9	246.7 1	300.7 1	56.79	27.2 9	0.8 9	43.97	7.89	1.2 9	7.3 0
IV	4.9 9	0.1 7	1.7 2	71.44	342.0 1	619.8 2	41.81	25.0 6	0.6 4	85.70	15.2 1	1.7 8	5.9 2
V	4.6 8	0.1 4	2.1 6	33.84	102.8 7	336.0 0	60.32	28.4 9	0.9 9	134.7 3	24.8 9	3.1 3	4.7 5
VI	5.4 0	0.1 6	1.7 1	56.91	242.0 5	497.7 9	54.18	33.6 6	0.6 8	71.39	23.3 5	2.8 1	2.2 8
VII	5.2 2	0.1 5	1.5 2	58.70	217.0 7	470.0 1	57.12	26.7 7	0.8 0	70.42	21.2 8	2.5 7	2.0 2

340 Among the average and complete linkage methods under Euclidean and squared Euclidean  
 341 distances, clustering based on average linkage method provided better clusters as the DB index was  
 342 low. DB index of complete linkage under squared Euclidean was 0.427 which was higher than that of  
 343 average linkage (0.412). Average linkage method had small DB index (0.306) under Euclidean  
 344 distance compared to squared Euclidean.

#### 345 4. CONCLUSION

346 Soils exhibit a high degree of heterogeneity with respect to place and climate. Spatial  
 347 distribution of soil nutrients also vary with locations [21]. Fertilizer recommendations are given based  
 348 on the soil test results of the locations and if the locations having similar soil characteristics can be  
 349 grouped based on soil fertility status, recommendations could be given easily [22].

350 Clustering of panchayats in Ernakulam were carried out using hierarchical clustering and *K*- means  
 351 clustering. Different distance measures like euclidean distance, squared Euclidean, Chebychev  
 352 distance and Mahalanobis  $D^2$  along with different linkage methods – single, complete and average  
 353 linkage were used for clustering the panchayats based on soil fertility status. . Davies- Bouldin index  
 354 was determined for each clustering procedure and Mahalanobis  $D^2$  was found to be the best  
 355 clustering method followed by euclidean distance with average linkage method as the D-B index was  
 356 small for Mahalanobis  $D^2$ .

357 Clustering of panchayats in Ernakulam based on Mahalanobis  $D^2$  resulted in seven clusters.  
 358 Keerampara and Manjapra remained as individual clusters. Keerampara soils were strongly acidic  
 359 with EC value 0.06 dS m<sup>-1</sup>. Available Ca was observed to be comparatively lowest (274 mg kg<sup>-1</sup>)  
 360 while Fe content was the highest (158 mg kg<sup>-1</sup>). Available Mg was very low in Manjapra (19 mg kg<sup>-1</sup>)  
 361 compared to other panchayats. On the contrary, S content was more in this panchayat (57 mg kg<sup>-1</sup>).  
 362 Nedumbassery and Piravom belonged to the same cluster due to the similarities in their soil  
 363 characteristics. There were same EC (0.10 dS m<sup>-1</sup>), OC ranged from 1.2 to 1.5 per cent, K ranged  
 364 from 161 to 168 kg ha<sup>-1</sup> and S ranged from 20 to 21 mg kg<sup>-1</sup> in these two panchayats. Kakkad,  
 365 Kalady and Vengoor belonged to the same cluster as EC (0.15-0.19 dS m<sup>-1</sup>), OC (2-2.4%) and Mg  
 366 (71-73 mg kg<sup>-1</sup>) had approximately same values.

#### 367 ACKNOWLEDGEMENT

368 This research was supported by Kerala Agricultural University (KAU) and we thank our colleagues  
 369 from KAU who provided insight and expertise that greatly assisted the research. We would also like to  
 370 show our gratitude to all those who gave assistance and comments for completing this research.

## 371 REFERENCES

- 372 1. Das DK. Introductory Soil Science (Indian reprint, 2016).Kalyani Publishers, New Delhi;  
373 1996.
- 374 2. Bernstein L. Effects of salinity and sodicity on plant growth. Annu. Rev. Phytopathol. 1975;  
375 13:295-312.
- 376 3. Das MN, Giri NC. Design and analysis of experiments. 3rd ed. New age International  
377 Publishers, New Delhi; 1979.
- 378 4. Dawes L, Goonetilleke A. Using multivariate analysis to predict the behavior of soils under  
379 effluent irrigation. Water Air Soil Pollut. 2006; 172 (1-4): 109-127.
- 380 5. Bharadwaj C, Satyavathi CT, Subramanyam D. Evaluation of different classificatory analysis  
381 methods in some rice (*Oryza sativa*) collections. Indian J. Agric. Sci. 2001; 71(2): 123-125.
- 382 6. Temp M, Filzmoser P, Reimann C. Cluster analysis applied to regional geochemical data:  
383 Problems and possibilities. [on-line]. 2006. Available:  
384 <https://www.sciencedirect.com/science/article/pii/S088329270800125X>.
- 385 7. Altdorff D, Dietrich P. Cluster analysis of geophysical field data: an approach for reasonable  
386 partitioning of sites. 2010. In: 19th World Congress of Soil Science on 'Soil Solutions for a  
387 Changing World' 1 – 6 August 2010, Brisbane, Australia [On-line].
- 388
- 394 8. Ye R, Wright AL. Multivariate analysis of chemical and microbial properties in histosols as  
395 influenced by land-use types. Soil Till. Res. 2010; 100:94-100.
- 396 9. Morrison JM, Goldhaber MB, Ellefsen KJ, Mills CT. Cluster analysis of a regional-scale soil  
397 geochemical dataset in northern California. J. Int. Assoc. Geochem. 2011; 26: S105-S107.
- 398 10. Legany C, Juhasz S, Babos B.. Cluster validity measurement techniques. Proceedings of the  
399 5th WSEAS Int. Conf. on Artificial Intelligence, Knowledge Engineering and Data Bases,  
400 Madrid, Spain, February 15-17, 2006; 388-393.
- 401 11. Ahmed F, Fakhruddin ANM, Imam MDT, Khan N, Khan TA, Rahman MM, *et al.* Spatial  
402 distribution and source identification of heavy metal pollution in roadside surface soil: a study  
403 of Dhaka Aricha highway, Bangladesh. Ecol. Process. 2016; 5(2): 1-16.
- 404 12. Bansod BS, Pandey OP, Rajesh NL. Analysis and delineation of spatial variability using geo-  
405 sensed apparent electrical conductivity and clustering techniques. Int. J. Agric. Biol. 2012;  
406 14: 481–491.
- 407 13. Granato D, Santosa JS, Eschera GB, Ferreira BL, Maggio RM. Use of principal component  
408 analysis (PCA) and hierarchical cluster analysis (HCA) for multivariate association between  
409 bioactive compounds and functional properties in foods: A critical perspective. Trends Food  
410 Sci. Technol. 2018; 72: 83–90.
- 411 14. Mahalanobis PC. On the generalized distance in statistics. Proceedings of the National  
412 Institute of sciences (Calcutta), 1936; 2: 49-55.
- 413 15. de Araujoneto JR, Gomes FEF, de Q Palacio HA, da Silva EB, Brasil PP. Similarity of soils  
414 with regard to salinity in the perennial valley of trussu river, Ceará. Irriga, Botucatu. 2016;  
415 21(2): 327-341.
- 416 16. Ellefsen KJ, Smith DB. Manual hierarchial clustering of regional geochemical data using a  
417 Bayesian finite mixture model. J. Int. Assoc. Geochem. 2016; 75: 200-210.
- 418 17. Guler C, Thyne GD, McCray JE, Turner AK. Evaluation of graphical and multivariate statistical  
419 methods for classification of water chemistry data. Hydrogeol. J. 2002; 10:455–474.
- 420 18. Bakry AB, Ibrahim OM, Abd El-Fattah Elewa T, El-Karamany MF. Performance assessment  
421 of some flax (*Linum usitatissimum* L.) varieties using cluster analysis under sandy soil  
422 conditions. Agric. Sci. 2014; 5: 677-686.
- 423 19. Ansari Z, Babu AV, Azeem MF, Ahmed W. Quantitative evaluation of performance and  
424 validity indices for clustering the web navigational sessions. World Comput. Sci. Inf. Technol.  
425 J. 2011; 1(5): 217-226.

- 426 20. GoK (Government of Kerala). Fertility of Soils of Kerala. State planning board,  
427 Thiruvananthapuram; 2013.
- 428 21. Al-Farhud A, Al-Sewailem M, Usman ARA. Status of selenium and trace elements in some  
429 arid soils cultivated with forage plants: A case study from Saudi Arabia. *Int. J. Agric. Biol.*  
430 2017; 19: 85–92.
- 431 22. KAU (Kerala Agricultural University). Package of Practices Recommendations: Crops (15th  
432 Ed.). Kerala Agricultural University, Thrissur; 2016.

433

434

435

436

437