

An Exploration of the DBSCAN Algorithm

Stat 666 Term Project

Matthew Morgan, Jace Ritchie

Abstract

Density-based spatial clustering of applications with noise (DBSCAN) is a popular clustering algorithm. In this report, simulations prove the efficacy of DBSCAN over k -means as dimensionality scales over data generated with non-overlapping means, as well as the shared inefficacy of both methods when the data mean vectors overlap significantly. DBSCAN is also used to cluster a dataset of player/course performances from the 2014 season of the Professional Golf Association.

1 Introduction

In the field of unsupervised learning, algorithms are used to analyze and cluster unlabeled datasets. Each algorithm uses a specific clustering process in an effort to discover patterns and/or groupings in the data that might not be as evident to the human eye. The patterns and/or groupings in the data found by these algorithms are used in a wide range of applications.

One clustering algorithm that has become more popular is the density-based spatial clustering of applications with noise (DBSCAN) algorithm. DBSCAN is a popular clustering algorithm that can discover clusters of arbitrary shape and size while also identifying noise points. Due to its robustness and scalability, DBSCAN is often chosen for clustering tasks in various domains such as image processing, natural language processing, and social network analysis. This report aims to provide an in-depth exploration of the DBSCAN algorithm, including its core concepts, implementation, advantages, and limitations.

Specifically, this exploration of the DBSCAN algorithm includes a simulation study which compares the performance of the DBSCAN algorithm to the commonly used k -means algorithm. The simulation study focuses on evaluating the two algorithms with respect to their multidimensional scaling ability as well as their ability to cluster data generated from various settings. Additional exploration of the DBSCAN algorithm includes using it with a real-world dataset from the 2014 season of the Professional Golf Association to demonstrate its effectiveness and applicability. The insights gained from this exploration of the DBSCAN algorithm can be useful to practitioners and researchers interested in clustering techniques and can serve as a guide for selecting appropriate clustering methods for different types of data.

2 Methodology

2.1 Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

DBSCAN is a clustering methodology based on the distance between points rather than a pre-allocated number of clusters as defined by Ester et al. (1996). In contrast to many popular clustering algorithms, instead of minimizing within-group variance and maximizing across-group variance with respect to a defined number of clusters, DBSCAN finds a network of points that are a given distance away from each other, ϵ . Clusters are defined when there is a group of points of at least a minimum size ($MinPts$) within ϵ of one of the points. From there, clusters are expanded by adding points within ϵ to the clusters. Thus, instead of finding a predetermined number of clusters, the number of clusters that DBSCAN finds is determined by the number of clusters found to be of at least size $MinPts$ with a density directly tied to ϵ . Furthermore, DBSCAN does not require that every point belong to a cluster. As a result, there may be some points that are simply labelled as “noise”.

Algorithm 1 DBSCAN

Object: Given a set of data points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ where $\mathbf{x}^{(i)} \in \mathbb{R}^n$, identify clusters with some minimum number of points that are a specified distance away from each other.

Procedure:

- 1: Set ϵ and *MinPts* to be the specified distance and minimum number of points necessary to make a cluster, respectively.
 - 2: Choose a random, unclustered point, p , that hasn't been selected yet. If such a point doesn't exist, go to Step 7.
 - 3: Identify all t points that are closer than ϵ to p .
 - 4: If p is a cluster or $t > \text{MinPts}$, then p and its t nearest neighbors form a cluster. Otherwise, go to Step 2.
 - 5: If $t = 0$, go to Step 2.
 - 6: Define p to be p and its t nearest neighbors. Go to Step 3.
 - 7: Define all unclustered points as noise.
-

As a rule of thumb, *MinPts* is generally set to twice the number of variables being considered (Sander et al., 1998). Additionally, ϵ is set to be the inflection point on a *MinPts* – 1 nearest neighbors distance plot, as shown in Figure 1 (Schubert et al., 2017). These heuristic methods will be used throughout the remainder of the report to find useful parameters for DBSCAN.

2.2 k -means

The other clustering algorithm considered in this project was k -means clustering (Renchner and Christensen, 2012). This method was considered as it is one of the most commonly used procedures in unsupervised learning. Given the prevalent use of the k -means clustering method, it was used as a baseline to which the DBSCAN clustering method could be compared. The details of the k -means algorithm are the following:

Algorithm 2 k -means

Object: Given a set of data points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ where $\mathbf{x}^{(i)} \in \mathbb{R}^n$, predict k centroids and a label $c^{(i)}$ for each data point.

Procedure:

- 1: Initialize cluster centroids $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$ randomly.
- 2: Assign each data point, $i \in 1, \dots, m$, to closest cluster centroid (using Euclidean distance).

$$c^{(i)} = \arg \min_j \|\mathbf{x}^{(i)} - \mu_j\|^2$$

- Recalculate each of the cluster centroids, $j \in 1, \dots, k$.

$$\mu_j = \frac{\sum_{i=1}^m 1\{c^{(i)} = j\} \mathbf{x}^{(i)}}{\sum_{i=1}^m 1\{c^{(i)} = j\}}$$

- 3: Repeat Step 2 until no more reassignments occur.
-

One interesting aspect of the k -means algorithm is that the number of clusters into which data points will be labeled must be chosen *a priori*. Therefore, there must be some sort of intuition beforehand about the underlying grouping of the data points when implementing the k -means algorithm.

3 Simulation Study

3.1 Setup

As a means of evaluating the performance of the DBSCAN algorithm and the k -means algorithm, a simulation study was conducted where the “true” underlying grouping of objects was set *a priori*. The simulation was conducted using 10,000 iterations. The data for the four groups were generated from four multivariate normal distributions, each with 100 data points, equal covariance matrices

$$\Sigma_{p_s \times p_s} = \begin{bmatrix} 5 & .5 & \dots & .5 \\ .5 & 5 & \dots & .5 \\ \vdots & \vdots & & \vdots \\ .5 & .5 & \dots & 5 \end{bmatrix},$$

and distinct $p_s \times 1$ mean vectors μ_{1s} , μ_{2s} , μ_{3s} , and μ_{4s} for $s \in 1, \dots, 4$. Additionally, the data was standardized after being generated.

- The first set of mean vectors used to generate the data were of length $p_1 = 4$ with:
 - $\mu_{11} = (5, 6, 7, 8)$, a baseline mean vector
 - $\mu_{21} = (-5, -6, -7, -8)$, a negatively collinear mean vector
 - $\mu_{31} = (15, 18, 21, 24)$, a positively collinear mean vector
 - $\mu_{41} = (0, 0, 0, 0)$, a non-collinear mean vector
- The second set of mean vectors used to generate the data were of length $p_2 = 4$ with:
 - $\mu_{12} = (5, 6, 7, 8)$, a baseline mean vector
 - $\mu_{22} = (3, 4, 5, 6)$, an overlapping collinear mean vector
 - $\mu_{32} = (7, 8, 9, 10)$, an overlapping collinear mean vector
 - $\mu_{42} = (0, 0, 0, 0)$, a non-overlapping mean vector
- The third set of mean vectors used to generate the data were of length $p_3 = 8$ with:
 - $\mu_{13} = (5, 6, 7, 8, 9, 10, 11, 12)$, a baseline mean vector
 - $\mu_{23} = (-5, -6, -7, -8, -9, -10, -11, -12)$, a negatively collinear mean vector
 - $\mu_{33} = (15, 18, 21, 24, 27, 30, 33, 36)$, a positively collinear mean vector
 - $\mu_{43} = (0, 0, 0, 0, 0, 0, 0, 0)$, a non-collinear mean vector
- The fourth set of mean vectors used to generate the data were of length $p_4 = 12$ with:
 - $\mu_{14} = (5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16)$, a baseline mean vector
 - $\mu_{24} = (-5, -6, -7, -8, -9, -10, -11, -12, -13, -14, -15, -16)$, a negatively collinear mean vector
 - $\mu_{34} = (15, 18, 21, 24, 27, 30, 33, 36, 39, 42, 45, 48)$, a positively collinear mean vector
 - $\mu_{44} = (0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0)$, a non-collinear mean vector

The simulation study was conducted using mean vectors of increasing length to examine the multidimensional scaling ability of the two algorithms. Then, to assess the ability of the two clustering algorithms to group data generated from different mean vector settings, two sets of mean vectors of length 4 were used to generate the data with one set being non-overlapping and the other set overlapping.

For each simulation study setting s , values of $MinPts$ and ϵ were chosen for the DBSCAN algorithm according to the heuristics given in Section 2.1. The k -nearest neighbor distance plots used to select the values of ϵ for each simulation study setting are shown in Figure 1.

Within each iteration of the simulation study, after the data was generated and standardized, the DBSCAN algorithm was then implemented to cluster the data. Then, due to the fact that the DBSCAN algorithm automatically identifies the number of clusters in the data without prior knowledge, it was first checked whether the DBSCAN algorithm found the correct number of “true” underlying clusters.

Next, if the DBSCAN algorithm had found the correct number of “true” underlying clusters, the percent of the data points left unclustered was recorded. Then, on the set of data points that the DBSCAN algorithm

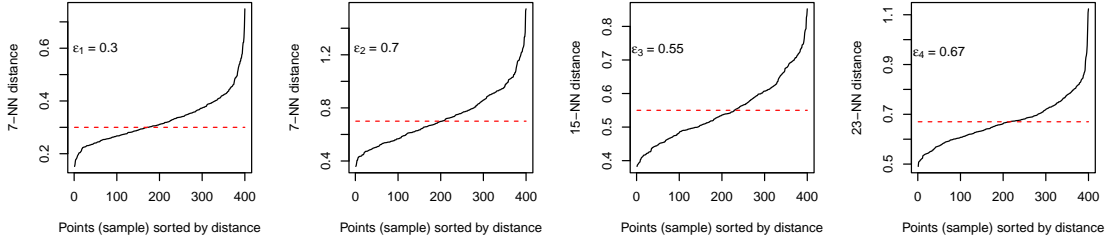


Figure 1: k -nearest neighbor distance plots for $s \in \{1, 2, 3, 4\}$ and their respective values of ϵ .

had clustered, k -means was utilized to cluster that same set of data points into 4 clusters as a comparison. Finally, the accuracy of group identification, as measured by the proportion of data points assigned to the correct group was recorded to evaluate the performance of the two clustering algorithms.

3.2 Results

Table 1: Simulation Study Results

s	Found Groups	Percent Noise	DBSCAN Accuracy	k -means Accuracy
1	0.9113	0.2067	0.8991	0.7228
2	0.0037	0.0008	0.0012	0.0017
3	1.0000	0.0523	1.0000	0.8802
4	1.0000	0.0225	1.0000	0.7892

An examination of Table 1 reveals that in the settings that used non-overlapping mean vectors ($s \in \{1, 3, 4\}$), there were four key results. The first result to note is that the DBSCAN algorithm was able to find the “true” underlying number of groups a significant majority or all of the time. The second result to note is that the DBSCAN algorithm classified less data points as noise as the dimensionality of the mean vectors increased. This phenomenon makes sense as the mean vectors were increasing in complexity as the dimensionality increased, which allowed for greater values of ϵ . This in turn caused the DBSCAN algorithm to classify fewer data points as noise. The third result to note is that the DBSCAN algorithm improved in its classification accuracy of the data points as the dimensionality of the mean vectors increased. This trend did not hold for the k -means algorithm. This behavior in the DBSCAN algorithm is most likely a result of the fact that the mean vectors were less likely to overlap as the dimensionality increased. The final key result to mention is that the DBSCAN algorithm more accurately classified the data points than the k -means algorithm did in all three settings.

A further examination of Table 1 reveals that in the setting that used heavily overlapping mean vectors ($s = 2$), both the DBSCAN algorithm and the k -means algorithm performed poorly. Specifically, the DBSCAN algorithm rarely found the correct number of “true” underlying groups and was widely inaccurate in its classification of the data points. It should also be noted that the k -means algorithm was just as inaccurate in its classification of the data points.

Figure 2 provides a visual representation of the clusters found by the DBSCAN algorithm in the settings of non-overlapping and overlapping mean vectors ($s \in \{1, 2\}$). An examination of Figure 2 further highlights the inefficiency of the DBSCAN algorithm to find distinct clusters in the setting where overlapping mean vectors were used to generate the data as it cannot separate out the groups. Given this inefficiency of the DBSCAN algorithm, it gives the indication that the DBSCAN algorithm relies on separation to be able to identify distinct clusters.

Overall, the simulation study conducted to evaluate the performances of the DBSCAN algorithm and the k -means algorithm revealed that the DBSCAN algorithm is robust to multidimensional scaling and excels

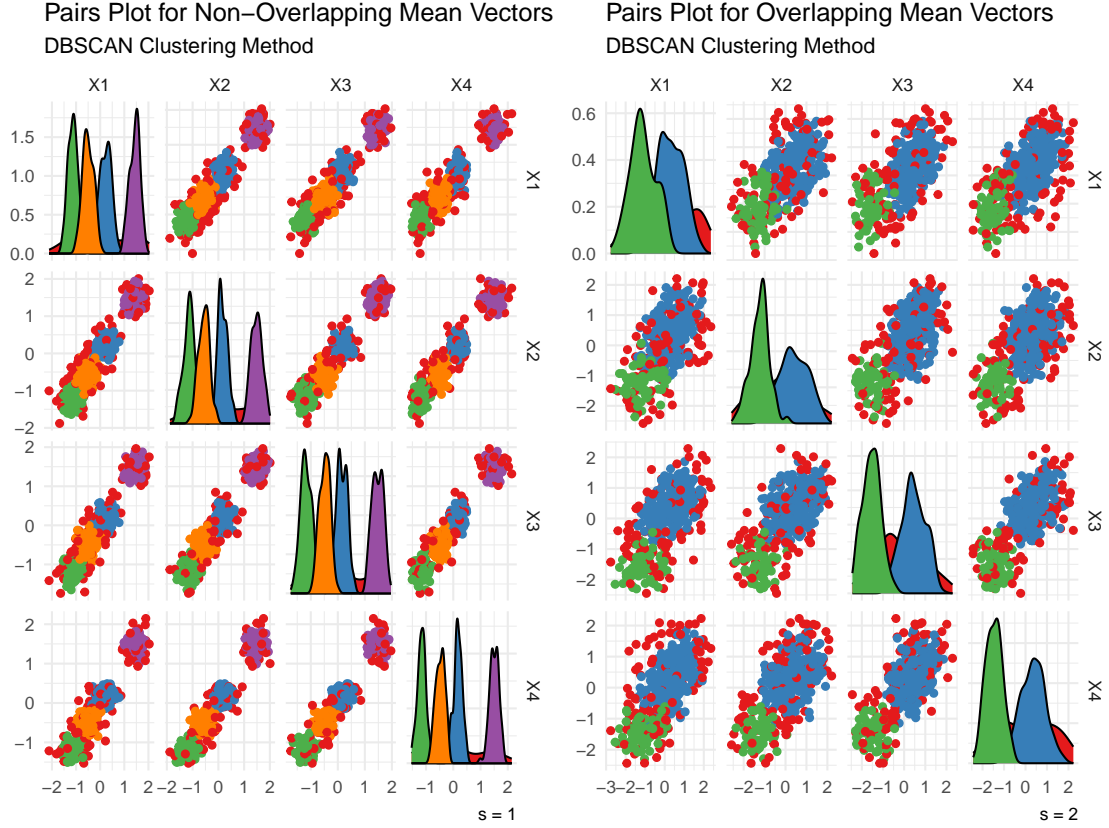


Figure 2: Pairs plots for one iteration of the simulation study for the first two settings, $s \in \{1, 2\}$. The red points represent the data points that the DBSCAN algorithm classified as noise.

at clustering when sufficient separation exists. In the settings where mean vectors with a non-overlapping structure were used to generate the data, the DBSCAN algorithm consistently outperformed the k -means algorithm. Finally, the simulation study revealed a shared weakness between the DBSCAN and k -means algorithms that both methods perform poorly when the data was generated from heavily overlapping mean vectors.

4 Application

4.1 Dataset

The real-world dataset that was used with the DBSCAN algorithm comes from the 2014 season of the Professional Golf Association (PGA). It included information such as the player name, the course being played on, the slope of the course, the players' strokes gained (SG) in several categories, and the score on a common par of 75 (CP). There were 1,086 observations, with 22 different players and 18 different courses. Every course played by a given athlete was played four times by that athlete. The columns used for clustering were SG while putting, SG from tee to green, SG on approach, SG around the green, CP, course slope, player's average driving distances, and the number of hits in regulation onto the green.

4.2 Results

After finding the appropriate values of $MinPts$ and ϵ using the heuristic methods of Section 2.1 to be 16 and 1.4 respectively, the DBSCAN algorithm was used to cluster the data. Using DBSCAN on the data after standardization, 16 clusters were identified with a median group size of 41.5 observations.

Table 2: Group attributes of clusters identified by DBSCAN for the golf data

Cluster	Noise	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Observations	182	49	40	44	343	43	34	26	16	40	35	24	26	45	47	48	44
Players	19	1	1	1	9	1	1	1	1	1	1	1	2	1	1	2	1
Courses	18	13	11	13	18	11	11	9	5	13	12	7	9	12	14	13	12

After examining the contents of each cluster in Table 2, it was evident that the majority of the clusters (with the exception of cluster 4) were composed of a single player across several courses.

It is significant that three or four observations of a course were included in the same cluster for the same player, as this indicated that the player performed similarly on each of the courses. Overall these results were encouraging, since they implied that the DBSCAN algorithm found defined groups not included in the provided data, such as course and player. More specifically, with cluster 4, given that it contained roughly a third of the player/course combinations which were fairly similar, this cluster could be interpreted as an "average player on an average course" group. Finally, the 182 noise points identified by the DBSCAN algorithm can be thought of as outliers representing an instance where a player's performance on a particular course was distinct from all other performances in the dataset in some way.

5 Discussion

As shown in the simulations, the DBSCAN algorithm is quite competent at identifying collinear, but distinct groups even as the dimensionality of the data increases. It has some failings when identifying groups with overlapping means, but seemed to have similar failings to the k -means algorithm in those situations. Additionally, when applied to a real-world dataset, it seemed to identify known groups and find relationships that implied similarity between group members. Therefore, there seems to be practical use for this algorithm in both simulated and real-world scenarios.

After conducting this analysis, if all points are desired to be clustered, one proposal would be to use the centroids of the clusters identified by the DBSCAN algorithm as the initial cluster centroids in the k -means algorithm.

While this analysis explored some aspects of the DBSCAN algorithm, it is suggested that there be further exploration to continue to evaluate the robustness of this algorithm. Given that the simulation study in this analysis generated "true" underlying groups of equal size from normal data with identical covariance matrices, future work should include simulation studies where the data was generated from other settings. Generating data from non-normal distributions, generating data with groups of unequal size, and generating data with dissimilar covariance matrices are all data generation methods that should be considered in future simulation studies to further explore the abilities of the DBSCAN algorithm.

References

- Ester, M., Kriegel, H.-P., Sander, J., Xu, X., et al. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pages 226–231.
- Rencher, A. C. and Christensen, W. F. (2012). *Cluster Analysis*, page 532–539. Wiley.
- Sander, J., Ester, M., Kriegel, H.-P., and Xu, X. (1998). Density-based clustering in spatial databases: The algorithm gbscan and its applications. *Data mining and knowledge discovery*, 2:169–194.
- Schubert, E., Sander, J., Ester, M., Kriegel, H. P., and Xu, X. (2017). Dbscan revisited, revisited: Why and how you should (still) use dbscan. *ACM Trans. Database Syst.*, 42(3).

Appendix A: Analysis Code

```
1 # Reading in necessary packages
2 library(dbscan)
3 library(dplyr)
4 library(ggplot2)
5 library(grid)
6 library(gridExtra)
7
8 #####
9 ## Simulation Study ##
10 #####
11
12 #####
13 ## Setup ##
14 #####
15
16 truth <- c(rep(1, 100), rep(2, 100), rep(3, 100), rep(4, 100))
17
18 # Collinear mean vectors
19 mean_vecs1 <- list(
20   c(5:8),
21   c(-5:-8),
22   c(seq(from = 15, to = 24, by = 3)),
23   c(rep(0, 4))
24 )
25
26 # Overlapping mean vectors
27 mean_vecs2 <- list(
28   c(5:8),
29   c(3:6),
30   c(7:10),
31   c(rep(0, 4))
32 )
33
34 mean_vecs3 <- list(
35   c(5:12),
36   c(-5:-12),
37   c(seq(from = 15, to = 36, by = 3)),
38   c(rep(0, 8))
39 )
40
41 mean_vecs4 <- list(
42   c(5:16),
43   c(-5:-16),
44   c(seq(from = 15, to = 48, by = 3)),
45   c(rep(0, 12))
46 )
47
48 # Covariance matrices for each setting of mean vectors
49 cov_mat1 <- matrix(.5, nrow = 4, ncol = 4)
50 cov_mat1[row(cov_mat1) == col(cov_mat1)] = 5
51
52 cov_mat2 <- matrix(.5, nrow = 4, ncol = 4)
53 cov_mat2[row(cov_mat2) == col(cov_mat2)] = 5
54
55 cov_mat3 <- matrix(.5, nrow = 8, ncol = 8)
56 cov_mat3[row(cov_mat3) == col(cov_mat3)] = 5
57
58 cov_mat4 <- matrix(.5, nrow = 12, ncol = 12)
59 cov_mat4[row(cov_mat4) == col(cov_mat4)] = 5
60
61 n_vec <- rep(100, 4)
62
63 # par(mfrow = c(1, 4), oma = c(0, 0, 2, 0))
64
65 set.seed(16)
```



```

66 #####
67 # Setting 1 #
68 #####
69
70
71 # Generate data
72 X_g1 <- mvtnorm::rmvnorm(n_vec[1], mean_vecs1[[1]], cov_mat1)
73 X_g2 <- mvtnorm::rmvnorm(n_vec[2], mean_vecs1[[2]], cov_mat1)
74 X_g3 <- mvtnorm::rmvnorm(n_vec[3], mean_vecs1[[3]], cov_mat1)
75 X_g4 <- mvtnorm::rmvnorm(n_vec[4], mean_vecs1[[4]], cov_mat1)
76
77 # Standardize data
78 X <- rbind(X_g1, X_g2, X_g3, X_g4)
79 X <- scale(X)
80
81 # Pick eps for DBSCAN
82 kNNdistplot(X, minPts = 8)
83 lines(x = 1:400, y = rep(.3, 400), lty = 'dashed', col = "red")
84 text(x = 50, y = .6, expression(paste(epsilon[1], " = 0.3")))
85
86 # Use dbscan to cluster data
87 db_res1 <- dbscan(X, eps = .3, minPts = 8)
88
89 pairs1 <- GGally::ggpairs(
90   data.frame(X),
91   columns = 1:4,
92   mapping = aes(color = as.factor(db_res1$cluster + 1L)),
93   upper = list(continuous = "points")
94 ) +
95   scale_color_brewer(palette = "Set1") +
96   scale_fill_brewer(palette = "Set1") +
97   theme_minimal() +
98   labs(
99     title = "Pairs Plot for Collinear Mean Vectors",
100     subtitle = "DBSCAN Clustering Method",
101     caption = expression(paste(s, " = 1"))
102   )
103
104 #####
105 # Setting 2 #
106 #####
107
108 # Generate data
109 X_g1 <- mvtnorm::rmvnorm(n_vec[1], mean_vecs2[[1]], cov_mat2)
110 X_g2 <- mvtnorm::rmvnorm(n_vec[2], mean_vecs2[[2]], cov_mat2)
111 X_g3 <- mvtnorm::rmvnorm(n_vec[3], mean_vecs2[[3]], cov_mat2)
112 X_g4 <- mvtnorm::rmvnorm(n_vec[4], mean_vecs2[[4]], cov_mat2)
113
114 # Standardize data
115 X <- rbind(X_g1, X_g2, X_g3, X_g4)
116 X <- scale(X)
117
118 # Pick eps for DBSCAN
119 kNNdistplot(X, minPts = 8)
120 lines(x = 1:400, y = rep(.7, 400), lty = 'dashed', col = "red")
121 text(x = 50, y = .9, expression(paste(epsilon[2], " = 0.7")))
122
123 # Use dbscan to cluster data
124 db_res2 <- dbscan(X, eps = .7, minPts = 8)
125
126 pairs2 <- GGally::ggpairs(
127   data.frame(X),
128   columns = 1:4,
129   mapping = aes(color = as.factor(db_res2$cluster + 1L)),
130   upper = list(continuous = "points")
131 ) +
132   scale_color_brewer(palette = "Set1") +
133   scale_fill_brewer(palette = "Set1") +

```

```

134 theme_minimal() +
135 labs(
136   title = "Pairs Plot for Overlapping Mean Vectors",
137   subtitle = "DBSCAN Clustering Method",
138   caption = expression(paste(s, " = 2"))
139 )
140
141 g1 <- grid.grabExpr(print(pairs1))
142 g2 <- grid.grabExpr(print(pairs2))
143 grid.arrange(g1, g2, widths = c(0.5, 0.5))
144
145 #####
146 # Setting 3 #
147 #####
148
149 # Generate data
150 X_g1 <- mvtnorm::rmvnorm(n_vec[1], mean_vecs3[[1]], cov_mat3)
151 X_g2 <- mvtnorm::rmvnorm(n_vec[2], mean_vecs3[[2]], cov_mat3)
152 X_g3 <- mvtnorm::rmvnorm(n_vec[3], mean_vecs3[[3]], cov_mat3)
153 X_g4 <- mvtnorm::rmvnorm(n_vec[4], mean_vecs3[[4]], cov_mat3)
154
155 # Standardize data
156 X <- rbind(X_g1, X_g2, X_g3, X_g4)
157 X <- scale(X)
158
159 # Pick eps for DBSCAN
160 kNNdistplot(X, minPts = 16)
161 lines(x = 1:400, y = rep(.55, 400), lty = 'dashed', col = "red")
162 text(x = 50, y = .8, expression(paste(epsilon[3], " = 0.55")))
163
164 #####
165 # Setting 4 #
166 #####
167
168 # Generate data
169 X_g1 <- mvtnorm::rmvnorm(n_vec[1], mean_vecs4[[1]], cov_mat4)
170 X_g2 <- mvtnorm::rmvnorm(n_vec[2], mean_vecs4[[2]], cov_mat4)
171 X_g3 <- mvtnorm::rmvnorm(n_vec[3], mean_vecs4[[3]], cov_mat4)
172 X_g4 <- mvtnorm::rmvnorm(n_vec[4], mean_vecs4[[4]], cov_mat4)
173
174 # Standardize data
175 X <- rbind(X_g1, X_g2, X_g3, X_g4)
176 X <- scale(X)
177
178 # Pick eps for DBSCAN
179 kNNdistplot(X, minPts = 24)
180 lines(x = 1:400, y = rep(.67, 400), lty = 'dashed', col = "red")
181 text(x = 50, y = .95, expression(paste(epsilon[4], " = 0.67")))
182
183 # Add a main title to the layout
184 mtext(
185   "Figure 1: kNN Dist Plots with Varying minPts",
186   outer = TRUE,
187   cex = 1.5,
188   line = 0,
189   font = 2
190 )
191
192 # Add a subtitle to the layout
193 mtext(
194   expression(paste("Red dashed line represents ", epsilon)),
195   outer = TRUE,
196   cex = 1,
197   line = -2
198 )
199
200 #make into a function so we can change by mean vectors
201 #do simulations three times, once for mean vectors of length 4,5, and 6 respectively

```

```

202 set.seed(16)
203 do_sim <- function(mean_vecs, n_vec, cov_mat, eps) {
204
205   sims <- 10000
206   n_groups <- 4
207   found_groups <- logical(sims)
208   db_acc <- numeric(sims)
209   kmean_acc <- numeric(sims)
210   db_rand <- numeric(sims)
211   kmean_rand <- numeric(sims)
212   perc_noise <- numeric(sims)
213
214   for (i in 1:sims) {
215
216     # Generate data
217     X_g1 <- mvtnorm::rmvnorm(n_vec[1], mean_vecs[[1]], cov_mat)
218     X_g2 <- mvtnorm::rmvnorm(n_vec[2], mean_vecs[[2]], cov_mat)
219     X_g3 <- mvtnorm::rmvnorm(n_vec[3], mean_vecs[[3]], cov_mat)
220     X_g4 <- mvtnorm::rmvnorm(n_vec[4], mean_vecs[[4]], cov_mat)
221     X <- rbind(X_g1, X_g2, X_g3, X_g4)
222
223     # Standardize data
224     X <- scale(X)
225
226     # Use dbscan to cluster data
227     db_res <- dbscan::dbscan(X, eps = eps, minPts = 2*length(mean_vecs[[1]]))$cluster
228
229     # Check if dbscan found 4 groups
230     if (max(db_res) != n_groups) {
231       found_groups[i] <- FALSE
232       next
233     }
234
235     # How many point did dbscan classify as noise?
236     perc_noise[i] <- sum(db_res == 0) / length(db_res)
237     found_groups[i] <- TRUE
238
239     # If 4 group were found by dbscan, compare it to k-means
240     k_res1 <- kmeans(X[db_res != 0,], n_groups)$cluster
241     k_res2 <- k_res1
242     u_k <- unique(k_res1)
243     k_res2[k_res1 == u_k[1]] <- 1
244     k_res2[k_res1 == u_k[2]] <- 2
245     k_res2[k_res1 == u_k[3]] <- 3
246     k_res2[k_res1 == u_k[4]] <- 4
247
248     # Check if groups were identified correctly bewteen the two methods
249     db_red <- db_res[db_res != 0]
250     db_acc[i] <- sum(db_red == truth[db_res != 0]) / length(db_red)
251     kmean_acc[i] <- sum(k_res2 == truth[db_res != 0]) / length(db_red)
252     # db_rand[i] <- fossil::adj.rand.index(db_red, truth[db_res != 0])
253     # kmean_rand[i] <- fossil::adj.rand.index(k_res2, truth[db_res != 0])
254
255   }
256
257   # Return metrics
258   return(
259     list(
260       perc_noise = perc_noise,
261       found_groups = found_groups,
262       db_acc = db_acc,
263       db_rand = db_rand,
264       kmean_acc = kmean_acc,
265       kmean_rand = kmean_rand)
266   )
267 }
268 }
269

```

```

270 sim1 <- do_sim(mean_vecs1, n_vec, cov_mat1, .3)
271
272 sim2 <- do_sim(mean_vecs2, n_vec, cov_mat2, .7)
273
274 sim3 <- do_sim(mean_vecs3, n_vec, cov_mat3, .55)
275
276 sim4 <- do_sim(mean_vecs4, n_vec, cov_mat4, .67)
277
278 data.frame(
279   setting = c(1:4),
280   founds_groups = c(mean(sim1[['found_groups']]), mean(sim2[['found_groups']]), mean(sim3[['found_groups']]), mean(sim4[['found_groups']])),
281   perc_noise = c(mean(sim1[['perc_noise']]), mean(sim2[['perc_noise']]), mean(sim3[['perc_noise']]), mean(sim4[['perc_noise']])),
282   dbscan_acc = c(mean(sim1[['db_acc']]), mean(sim2[['db_acc']]), mean(sim3[['db_acc']]), mean(sim4[['db_acc']])),
283   kmeans_acc = c(mean(sim1[['kmean_acc']]), mean(sim2[['kmean_acc']]), mean(sim3[['kmean_acc']]), mean(sim4[['kmean_acc']]))
284 ) %>%
285   knitr::kable(
286     format = "latex",
287     digits = 4,
288     col.names = c("$s$", "Found Groups", "Percent Noise", "DBSCAN Accuracy", "k-means Accuracy"),
289     booktabs = TRUE
290   )
291
292 #####
293 ## Application ##
294 #####
295
296 golf_dat <- read.csv('2014GolfMod1.csv')
297 X <- golf_dat[, c(11:14, 24:27)]
298
299 # Standardize data
300 X <- scale(X)
301
302 # Pick eps for DBSCAN
303 kNNdistplot(X, minPts = 16)
304
305 db_res <- dbscan(X, eps = 1.4, minPts = 16)
306 db_res
307
308 pairs(X, col = db_res$cluster + 1L)
309
310 golf_red <- golf_dat[db_res$cluster == 1,]
311
312 golf_dat$cluster <- db_res$cluster
313 golf_dat %>%
314   group_by(cluster) %>%
315   summarize(n_players = length(unique(Player.Name)), n_courses = length(unique(Course.Number))) %>%
316   print(n = 29)

```