

# Neural Networks Project 2022 - Question 2 Report

Kondylia Vergou, AM 4325

Panagiotis Vouzalis, AM 2653

Η συγγραφή του κωδικά έγινε σε περιβάλλον Windows

- Για να λειτουργήσουν οι εντολές που αναφέρονται παρακάτω σε περιβάλλον Linux πρέπει να αντικατασταθεί το < ; > με < : >
- Τα dependencies της εφαρμογής μας βρίσκονται στο φάκελο ./jars
- Ανάλυση με τη διαδρομή της πλατφόρμας που θα τρέξει η εφαρμογή, μπορεί να χρειαστεί να προστεθεί στο τέλος του αρχείου προς εκτέλεση το `java -cp ".;jars/*" main/Kmeans <dataset_name.txt> <number_of_clusters>`

Instructions for Windows Platforms

1. Open terminal inside the directory Neural Networks - Question 2
  2. Compile with: `javac -cp ".;jars/*" main/*.java`
  3. Create a new dataset with: `java -cp ".;jars/*" main/DatasetGenerator` - This is not necessary
  4. Run java application with: `java -cp ".;jars/*" main/Kmeans <dataset_name.txt> <number_of_clusters>`
- example: `java -cp ".;jars/*" main/Kmeans dataset.txt 9 20`

Instructions for Linux Platforms

1. Open terminal inside the directory Neural Networks - Question 2
  2. Compile with: `javac -cp ".:jars/*" main/*.java`
  3. Create a new dataset with: `java -cp ".:jars/*" main/DatasetGenerator` - This is not necessary
  4. Run java application with: `java -cp ".:jars/*" main/Kmeans <dataset_name.txt> <number_of_clusters>`
- example: `java -cp ".:jars/*" main/Kmeans dataset.txt 9 20`

Για τις ανάγκες της δεύτερης άσκησης δημιουργήσαμε το dataset.txt (παρεχεται και το dataset.csv) σύμφωνα με τις απαιτήσεις του ερωτηματολ.Σ2.

Ο κώδικας για τη δημιουργία του dataset βρίσκεται στο αρχείο ./main/DatasetGenerator.java

## Code overview: ./main/Kmeans.java

- Μ ομάδα < - to M ορίζεται απο το χρήστη:

```
java -cp ".;jars/*" main/Kmeans.java <dataset_name.txt> <number_of_clusters> <number_of_runs>

144 public static void main( String[] args) {
145
146     numberOfClusters = Integer.parseInt(args[1]);
```

- Φορτωμα αρχείου παραδείγματος:

```
157 // use arguments to create a new dataset or use argument "dataset_name.txt" to import an existing .txt dataset
158 DatasetGenerator dataset = new DatasetGenerator(args[0]);
159
160 //DatasetGenerator dataset = new DatasetGenerator("new"); // use if you want to create a new dataset
161 //DatasetGenerator dataset = new DatasetGenerator("dataset_name.txt"); // use if you want to import an existing dataset
162
163 dataset = dataset.getDataset();
164 //dataset.exportDatasetCSV();
165 //dataset.exportDatasetTXT();
```

- Data initialization - τυχαία επιλογή αρχικών centroids:

```
121 // choose the center points of our clusters at random
122 private static void initializeData() {
123     ArrayList<Integer> indexesOfCenterPoints = new ArrayList<>();
124     int indexOfCenterPoint;
125
126     for (int i = 0; i < numberOfClusters; i++) {
127
128         ArrayList<Point> cluster = new ArrayList<>();
129         clusters.add(cluster);
130         indexOfCenterPoint = randomGenerator.nextInt(dataset.size());
131
132         // ensure we don't pick the same center point twice
133         while (indexesOfCenterPoints.contains(indexOfCenterPoint)) {
134
135             indexOfCenterPoint = randomGenerator.nextInt(dataset.size());
136
137         }
138         indexesOfCenterPoints.add(indexOfCenterPoint);
139         centerPoints.add(dataset.get(indexOfCenterPoint));
140     }
141 }
142 }
```

- Υπολογισμός του απόστασης ομαδοποίησης (ευκλείδια απόσταση) για κάθε point και τοποθέτηση του point στο σύνολο με τη μικρότερη απόσταση από το κέντρο της αντίστοιχης ομάδας:

```
184 //
185 // compute euclidean distance between a point and each cluster center, for every point
186 for (Point point : dataset) {
187     ArrayList<Point> distancesFromCenters = new ArrayList<>();
188
189     // compute euclidean distance between a point and each cluster center
190     for (Point centerPoint : centerPoints) {
191         float distanceFromCenter = (float) Math.sqrt(Math.pow((point.getX() - centerPoint.getX()), 2) + Math.pow((point.getY() - centerPoint.getY()), 2));
192         distancesFromCenters.add(distanceFromCenter);
193     }
194
195     // add the point to the cluster that minimizes the distance from its respective cluster center
196     clusters.add(distancesFromCenters.indexOf(Math.min(distancesFromCenters)))>0 ? point : null;
197     pointFromCenterDistances.add(distancesFromCenters.indexOf(Math.min(distancesFromCenters)));
198 }
199 }
```

- Υπολογισμός των νέων κέντρων ως προς το μέσο όρο των στοιχείων κάθε cluster:

```
213 // get the new cluster center with respect to the mean value of the cluster points, for every cluster
214 for (ArrayList<Point> cluster : clusters) {
215
216     int clusterIndex = clusters.indexOf(cluster);
217     //System.out.println("Sum of cluster " + clusterIndex + ": Center = 0 + centerPoints.get(clusterIndex).toString());
218     myWriter.write("cluster " + clusterIndex + ", Center = " + centerPoints.get(clusterIndex).toString() + "\n");
219
220     float sumX = 0;
221     float sumY = 0;
222
223     int displayLimit = 0;
224     for (Point point : cluster) {
225
226         //System.out.println(point.toString() + " ");
227         myWriter.write(point.toString() + " ");
228         sumX += point.getX();
229         sumY += point.getY();
230
231         displayLimit++;
232         if (displayLimit == 10) {
233             myWriter.write("\n");
234             displayLimit = 0;
235         }
236     }
237
238     //System.out.println();
239     //System.out.println("Sum of cluster " + clusterIndex + ": " + sumX);
240     //System.out.println("Sum of cluster " + clusterIndex + ": " + sumY);
241     //System.out.println("Average of cluster " + clusterIndex + ": " + sumX / cluster.size() + ", " + sumY / cluster.size());
242
243     myWriter.write("\n");
244     //myWriter.write("Sum of cluster " + clusterIndex + ": " + sumX + "\n");
245     //myWriter.write("Sum of cluster " + clusterIndex + ": " + sumY + "\n");
246     //myWriter.write("Average of cluster " + clusterIndex + ": " + sumX / cluster.size() + ", " + sumY / cluster.size() + "\n");
247
248     // set the new cluster center
249     Point newCenter = new Point(sumX / cluster.size(), sumY / cluster.size());
250     clusters.add(newCenter);
251     myWriter.write("cluster " + clusterIndex + ", New center = " + centerPoints.get(clusterIndex).toString());
252     myWriter.write("cluster " + clusterIndex + ", New center = " + centerPoints.get(clusterIndex).toString() + "\n");
253 }
```

- Ελέγχος τερματισμού βάσει της συνθήκης "Έχουν μεταβληθεί τα διανυσματα w, μεταξύ 2 επαναλήψεων?":

```
257 // are the w_i vectors different between 2 runs?
258 // if yes then start a new iteration
259 // if no then end the run
260 for (int k = 0; k < pointFromCenterDistances.size(); k++) {
261
262     for (int i = 0; i < pointFromCenterDistances.get(k).size(); i++) {
263
264         if (previousPointFromCenterDistances.get(k).get(i).equals(pointFromCenterDistances.get(k).get(i))) {
265             areVectorsDifferent = false;
266             break;
267         }
268     }
269
270     if (areVectorsDifferent) {
271         break;
272     }
273
274 }
275
276 // if the w_i vectors are different then prepare for a new iteration
277 if (areVectorsDifferent) {
278     previousPointFromCenterDistances = pointFromCenterDistances;
279 }
280
281 clusters_of_every_run.add(clusters);
282 centers_of_every_run.add(centerPoints);
283
284 // clear the clusters after every iteration
285 if (i != (numberOfRuns-1)) {
286     for (ArrayList<Point> cluster : clusters) {
287         cluster.clear();
288     }
289 }
290
291 // if the w_i vectors are different between 2 runs then start a new iteration, else end
292 while (areVectorsDifferent) {
293 }
```

- Υπολογισμός του συνολικού σφάλματος ομαδοποίησης (ευκλείδια απόσταση) για όλα τα clusters:

```
301 // compute the total dispersion by adding the dispersions of each cluster
302 float totalDispersion = 0;
303 for (int d = 0; d < clusters.size(); d++) {
304
305     float clusterDispersion = 0;
306     for (float clusterDistances : pointFromCenterDistances.get(d)) {
307         clusterDispersion += clusterDistances;
308     }
309
310     totalDispersion += clusterDispersion;
311 }
312
313 // ArrayList<float> totalDispersions holds the total cluster dispersion for each run
314 totalDispersions.add(totalDispersion);
315
316 //System.out.println("\nRun " + (i+1) + ": " + "dispersion = " + totalDispersion);
317 myWriter.write("\nRun " + (i+1) + ": " + "dispersion = " + totalDispersion + "\n");
```

- 20 runs και αποθήκευση της λύσης με το ελάχιστο σφάλμα ομαδοποίησης:

```
java -cp ".;jars/*" main/Kmeans.java <dataset_name.txt> <number_of_clusters> <number_of_runs>

324 // keep the run with the minimum dispersion
325 if (totalDispersion < MIN_DISPERSION) {
326
327     //System.out.println("\nRun " + (i+1) + ": " + "totalDispersion = " + MIN_DISPERSION + "\n");
328     MIN_DISPERSION = totalDispersion;
329     run_of_MIN_DISPERSION = (i+1);
330     clusters_MIN_DISPERSION = clusters_of_every_run.get(i);
331     centerPoints_MIN_DISPERSION = centers_of_every_run.get(i);
332 }
333 }
```

- Αποθήκευση των 9 κέντρων συντεταγμένων των centroids και του σφάλματος ομαδοποίησης στο αρχείο final\_centroids.dispersions (παράδειγμα για 9 clusters και 20 runs):

----- K means clustering example -----  
Number of clusters = 9  
Number of runs = 20

----- Final centroid coordinates -----  
Cluster 1 center = 1.0416359,0.986985  
Cluster 2 center = 1.6022335,0.24239336  
Cluster 3 center = 1.2923727,1.799663  
Cluster 4 center = 1.7235991,1.5450246  
Cluster 5 center = 0.64948756,1.8053868  
Cluster 6 center = 0.3978845,0.25267282  
Cluster 7 center = 0.39896592,1.3871365  
Cluster 8 center = 0.17588888,1.7561957  
Cluster 9 center = 1.7966483,1.8517332

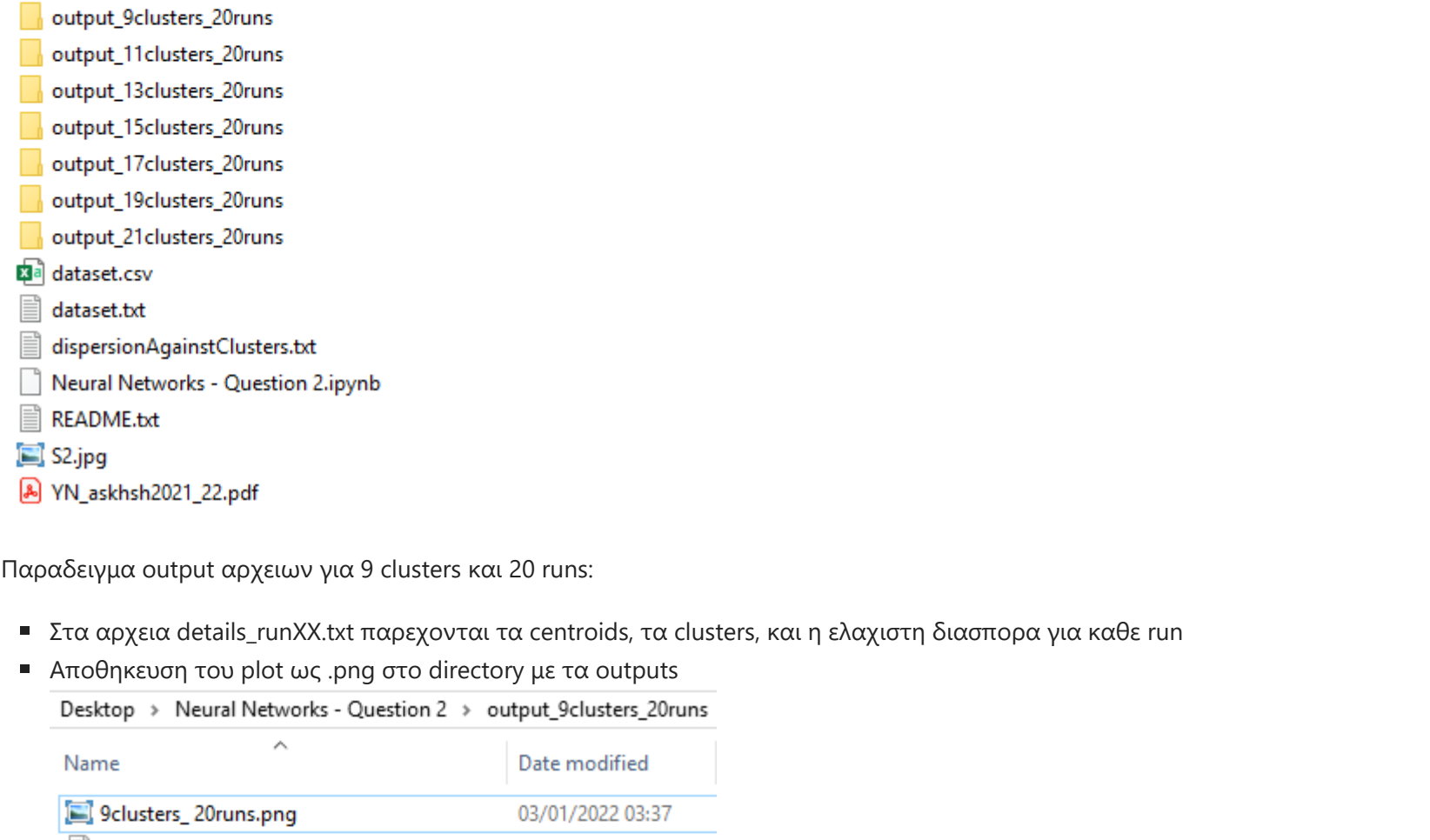
----- Dispersion across all runs -----  
Run 1 dispersion = 540.8662  
Run 2 dispersion = 384.42502  
Run 3 dispersion = 288.98016  
Run 4 dispersion = 261.02002  
Run 5 dispersion = 259.6171  
Run 6 dispersion = 259.16824  
Run 7 dispersion = 258.86148  
Run 8 dispersion = 258.8769  
Run 9 dispersion = 258.70828  
Run 10 dispersion = 258.57513  
Run 11 dispersion = 258.53192  
Run 12 dispersion = 258.57935  
Run 13 dispersion = 258.6531  
Run 14 dispersion = 258.70074  
Run 15 dispersion = 258.77325  
Run 16 dispersion = 258.51773  
Run 17 dispersion = 258.43878  
Run 18 dispersion = 257.83115  
Run 19 dispersion = 257.83108  
Run 20 dispersion = 257.83008

Minimum dispersion occurred on run 20 = 257.83008

- Δημιουργία plot μέσω της JFreeChart:

```
95 // panel generation
96 public Kmeans(String windowTitle, int numberOfClusters, int numberOfRuns, String dirString) {
97
98     super(windowTitle);
99
100     // create JFreeChart dataset
101     XYDataset dataset = createJFreeChartDataset();
102
103     // create scatter plot
104     String plotName = "K Means Clustering Example - " + numberOfClusters + " clusters & " + numberOfRuns + " runs";
105     JFreeChart chart = new JFreeChart(plotName, new ChartFactory.createScatterPlot(
106         plotName, "x axis", "y axis", dataset, PlotOrientation.VERTICAL, true, true, false));
107
108     // change background color
109     XYPlot plot = (XYPlot)chart.getPlot();
110     plot.setBackgroundPaint(new Color(255, 230, 200));
111
112     // create panel
113     ChartPanel panel = new ChartPanel(chart);
114     setContentPane(panel);
115
116     // export plot to png
117     try {
118         File imageFile = new File("./" + dirString + "/" + numberOfClusters + "clusters_" + numberOfRuns + "runs.png");
119         OutputStream out = new BufferedOutputStream(new FileOutputStream(imageFile));
120         BufferedImage image = chart.createBufferedImage(1000, 800);
121         ChartUtilities.writeBufferedImageAsPNG(out, image);
122     } catch (IOException e) {
123
124     }
125
126     System.out.println("An error occurred.");
127     e.printStackTrace();
128 }
129
130 // JFreeChart dataset creation
131 private XYDataset createJFreeChartDataset() {
132
133     XYSeriesCollection dataset = new XYSeriesCollection();
134     XYSeries centerSeries = new XYSeries("Centers");
135     for (Point point : centerPoints_MIN_DISPERSION) {
136         centerSeries.add(point.getX(), point.getY());
137     }
138
139     dataset.addSeries(centerSeries);
140
141     int i = 0;
142     for (ArrayList<Point> cluster : clusters_MIN_DISPERSION) {
143
144         XYSeries clusterSeries = new XYSeries("Cluster " + i);
145         for (Point point : cluster) {
146             clusterSeries.add(point.getX(), point.getY());
147         }
148         i++;
149         dataset.addSeries(clusterSeries);
150     }
151
152     return dataset;
153 }
```

- Δημιουργία γραφικού περιβάλλοντος για το plot:



- Αποθήκευση όλων των παραγόμενων αρχείων για κάθε στο root directory της εφαρμογής, κάθε run στο δικο του directory: Desktop > Neural Networks - Question 2

Name	
ipynb_checkpoints	
jars	
main	
output_3clusters_20runs	
output_7clusters_20runs	
output_9clusters_20runs	
output_13clusters_20runs	
output_17clusters_20runs	
output_21clusters_20runs	
dataset.csv	
dispersionAgainstClusters.txt	
Neural Networks - Question 2.ipynb	
README.txt	
52.jpg	
YN_askhsh2021_22.pdf	

- Παράδειγμα output αρχείων για 9 clusters και 20 runs:

- Στα αρχεία details\_runXX.txt παρέχονται τα centroids, τα clusters, και η ελάχιστη διασπορά για κάθε run
- Αποθήκευση του plot ως .png στο directory με τα outputs

Desktop > Neural Networks - Question 2 > output\_9clusters\_20runs

Clusters	Date modified
clusters_20runs.png	03/01/2022 09:37
details_run1.txt	03/01/2022 09:37
details_run2.txt	03/01/2022 09:37
details_run3.txt	03/01/2022 09:37
details_run4.txt	03/01/2022 09:37
details_run5.txt	03/01/2022 09:37
details_run6.txt	03/01/2022 09:37
details_run7.txt	03/01/2022 09:37
details_run8.txt	03/01/2022 09:37
details_run9.txt	03/01/2022 09:37
details_run10.txt	03/01/2022 09:37
details_run11.txt	03/01/2022 09:37
details_run12.txt	03/01/2022 09:37
details_run13.txt	03/01/2022 09:37
details_run14.txt	03/01/2022 09:37
details_run15.txt	03/01/2022 09:37
details_run16.txt	03/01/2022 09:37
details_run17.txt	03/01/2022 09:37
details_run18.txt	03/01/2022 09:37
details_run19.txt	03/01/2022 09:37
details_run20.txt	03/01/2022 09:37
final_centroids_dispersions.txt	03/01/2022 09:37
pd_minDispersion_clusters_20runs.txt	03/01/2022 09:37

## Μεταβολή του σφάλματος ομαδοποίησης αναλογα τον αριθμο των clusters

In [1]:

```
import pandas as pd
import matplotlib.pyplot as plt

df = pd.read_csv('dispersionAgainstClusters.txt')
df
```

Out[1]:

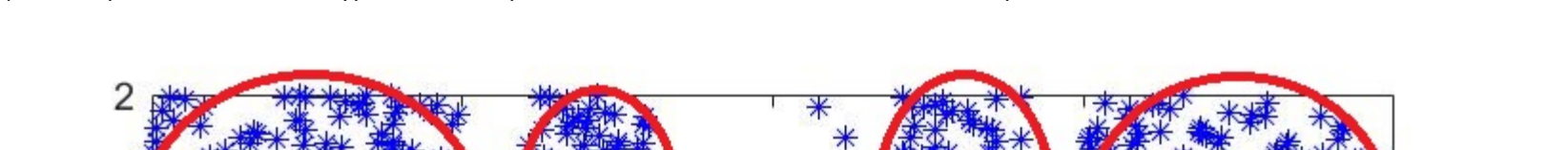
clusters	min-dispersion
0	3 578.31714
1	5 322.40290
2	7 271.04608
3	9 257.83008
4	11 226.77328
5	13 217.31060
6	15 181.79158
7	17 177.66852
8	19 164.44003
9	21 151.79976

In [2]:

```
# line plot of min dispersion against clusters
ax = df.plot(x='clusters', y='min-dispersion', figsize=(15, 10))
ax.set_title('Minimum dispersion against number of clusters')
ax.set_xlabel('number of clusters')
ax.set_ylabel('min dispersion')
```

Out[2]:

Text(0, 0.5, 'min dispersion')



## Μπορεί να χρησιμοποιηθεί το σφάλμα ομαδοποίησης για να εκτιμήσουμε τον πραγματικό αριθμο των cluster?

Σε γενικές γραμμές είναι δυνατόν να **εκτιμήσουμε** τον αριθμό των cluster εάν κανονίμε plot το σφάλμα ομαδοποίησης ως προς τις διαφοροποιήσεις του  $M$  (αριθμός cluster).

Συνήθως, όσο ο αριθμός  $M$  αυξάνεται η απόσταση μεταξύ των points σε κάθε cluster μειώνεται.

Όσοσο δεν μπορούμε να διαλέξουμε μια αυθαίρετα μεγάλη τιμή του  $M$ , διότι έτσι καταλήγουμε σε ένα από δυο πιθανά σεναρία:

1. Είτε θα υπάρχουν πολλά μικρά cluster με λίγα points.
2. Είτε θα υπάρχουν πολλά cluster με ένα μοναχικό point.

Either way, η ολή ιδέα της "ομαδοποίησης" αποδυναμώνεται όσο μεγαλώνει το  $M$ , για αυτο το λόγο θα πρέπει να διαλέξουμε μια "λογική" τιμή για το  $M$ .

Ένας τρόπος για να το κανονίμε αυτό είναι να χρησιμοποιήσουμε την ιδέα του "γονάτου" στο plot και να επιλέξουμε μια τιμή για την οποία το σφάλμα ομαδοποίησης μειώνεται δραματικά.

Στο παραπάνω plot παρατηρούμε ένα γονάτο στο διαστήμα  $M = 5$  έως  $M = 9$ .

Ειδικότερα, στην τιμή  $M = 5$  η τιμή του σφάλματος μειώνεται έντονα.

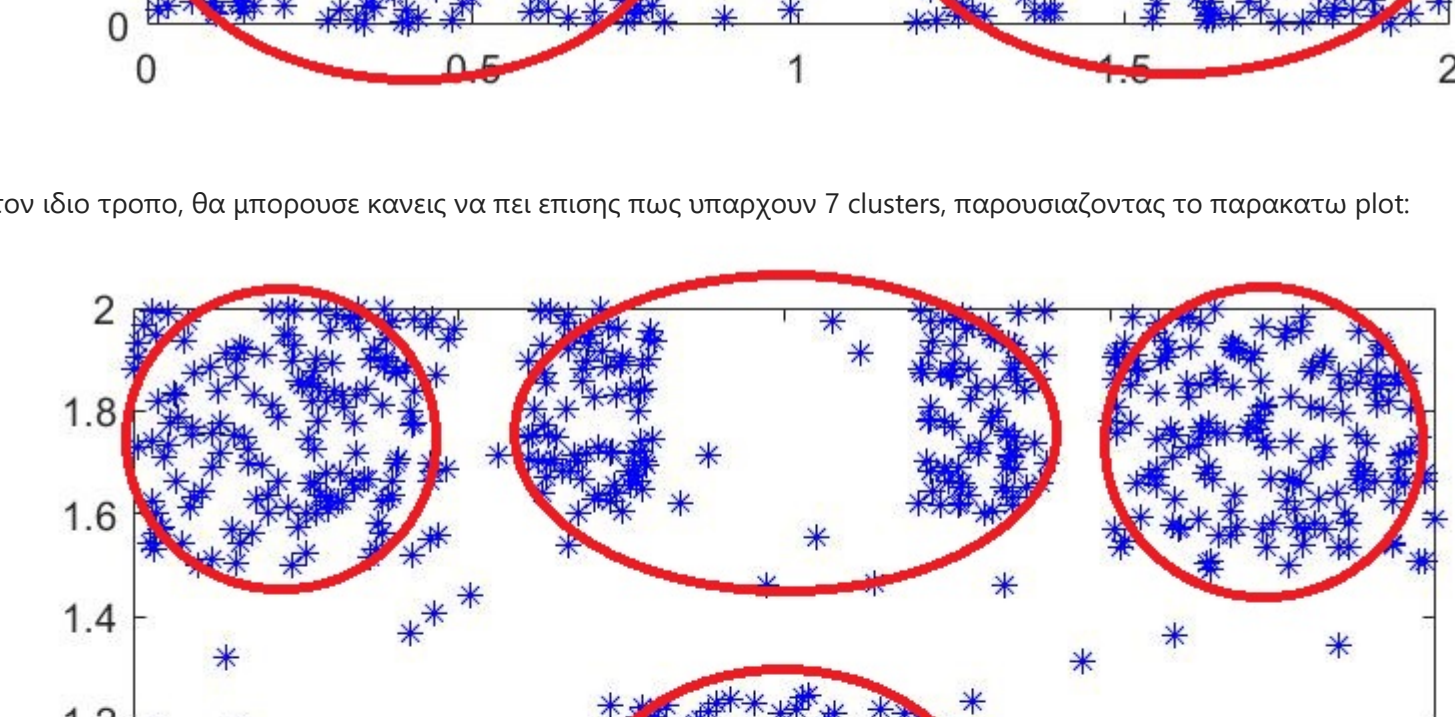
Στην τιμή  $M = 9$  μειώνεται λιγότερο έντονα.

Στην τιμή  $M = 7$  δεν μειώνεται τόσο έντονα όσο στις άλλες δυο.

Αρα αποσπάζουμε από αυτές τις τρεις τιμές  $M$  που είναι να χρησιμοποιήθει για την εκτίμηση του αριθμού των clusters, με τις πάντες  $M = 5$  είτε την τιμή  $M = 9$ .

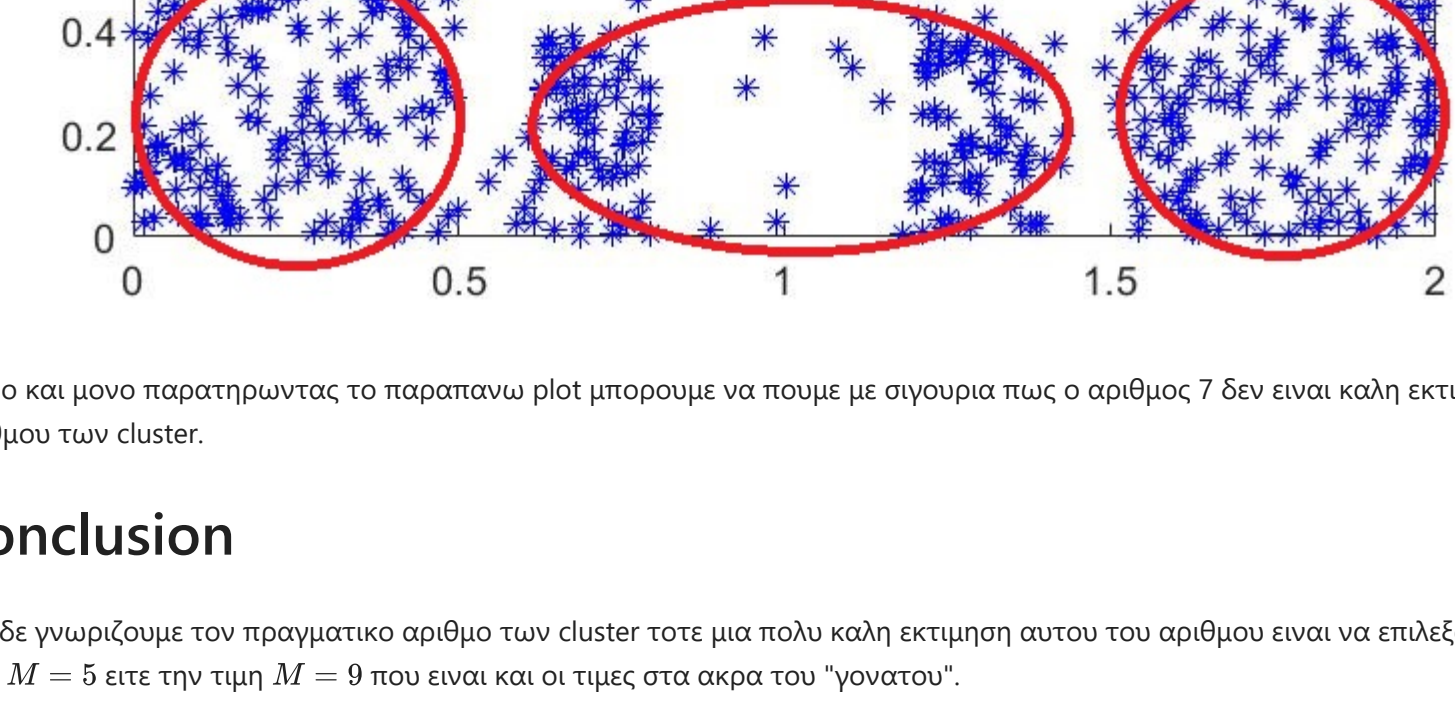
Εάν **δεν** είχαμε την προτερή γνώση πως ο πραγματικός αριθμός των clusters είναι 9, **τοτε πιθανότατα θα επιλέγαμε την τιμή 5** λόγω της έντονης πτώσης του σφάλματος.

Η εκκρίνιση μιας λείας πως ο πραγματικός αριθμός των clusters είναι 9. Ας το ερευνήσουμε:



Παρατηρούμε πως υπάρχουν πραγματικά 9 clusters στο dataset.

Επιχειρηματολογώντας για την ιδέα της υπαρκής 5 cluster, θα μπορούσε κανείς να πει πως υπάρχουν 5 cluster και να παρουσιάσει το παρακάτω plot:



Με τον ίδιο τρόπο, θα μπορούσε κανείς να πει επίσης πως υπάρχουν 7 clusters, παρουσιάζοντας το παρακάτω plot:



Μόνο και μόνο παρατηρώντας το παραπάνω plot μπορούμε να πούμε με αβυσσική πώς ο αριθμός 7 δεν είναι καλή εκτίμηση του αριθμού των cluster.

## Conclusion

Εάν δεν γνωρίζουμε τον πραγματικό αριθμό των cluster τότε μια πολύ καλή εκτίμηση αυτού του αριθμού είναι να επιλέξουμε είτε την τιμή  $M = 5$  είτε την τιμή  $M = 9$  που είναι και οι τιμές στα σκατά του "γονάτου".

Από τα παραπάνω plots βλέπουμε πως και οι δυο τιμές είναι "γονάτες", όμως από την εκκρίνιση γνωρίζουμε εκ των προτέρων πως η τιμή  $M = 9$  είναι η περισσότερο σωστή. Εν τέλει, όλα εξαρτώνται από το ποσο ελαστικό είμαστε όσον αφορά την "ομοιογένεια" των cluster που θέλουμε να φτιάξουμε (πιο αυθαίρετα μεγάλα ή → 5 clusters, μεγαλύτερη ομοιογένεια → 9 clusters).