

ΕΞΟΡΥΞΗ ΔΕΔΟΜΕΝΩΝ

Εργασία 3

ΠΑΝΤΕΛΕΗΜΩΝ ΠΡΩΙΟΣ

ice18390023

7ο Εξάμηνο

ice18390023@uniwa.gr

Τμήμα Τρίτης 9:00-13:00



ΠΑΝΕΠΙΣΤΗΜΙΟ ΔΥΤΙΚΗΣ ΑΤΤΙΚΗΣ
UNIVERSITY OF WEST ATTICA

Υπεύθυνοι καθηγητές

ΠΑΡΙΣ ΜΑΣΤΟΡΟΚΩΣΤΑΣ

Τμήμα Μηχανικών και Πληροφορικής Υπολογιστών
22 Ιανουαρίου 2022

Περιεχόμενα

1	Σύνολο δεδομένων του Enron	1
1.1	Δενδρόγραμμα για διαφορετικούς δεσμούς και δείκτες ομοιότητας	1
2	Διαμεριστική συσταδοποίηση με k-means	4
2.1	Διαχωριστική συσταδοποίηση για το σύνολο δεδομένων Iris	4
2.2	Διαχωριστική συσταδοποίηση για το σύνολο δεδομένων xV	15

Κώδικες

1.1	Script συσταδοποίησης των email του Enron	2
2.1	Script συσταδοποίησης με k-means για διάφορους μεθόδους στα δεδομένα Iris	14
2.2	Script συσταδοποίησης με k-means για διάφορους μεθόδους στα δεδομένα xV	17

Κατάλογος σχημάτων

1.1	Απλός δεσμός με ομοιότητα jaccard	1
1.2	Απλός δεσμός με ομοιότητα cosine	2
1.3	Μέσος δεσμός με ομοιότητα cosine	2
2.1	Διαχωρισμός συστάδων με euclidean για την τρίτη και τέταρτη στήλη	5
2.2	Διαχωρισμός συστάδων με cityblock για την τρίτη και τέταρτη στήλη	5
2.3	Διαχωρισμός συστάδων με cosine για την τρίτη και τέταρτη στήλη	6
2.4	Διαχωρισμός συστάδων με correlation για την τρίτη και τέταρτη στήλη	6
2.5	Διαχωρισμός συστάδων με euclidean για την πρώτη και δεύτερη στήλη	7
2.6	Διαχωρισμός συστάδων με cityblock για την πρώτη και δεύτερη στήλη	7
2.7	Διαχωρισμός συστάδων με cosine για την πρώτη και δεύτερη στήλη	8
2.8	Διαχωρισμός συστάδων με correlation για την πρώτη και δεύτερη στήλη	8
2.9	Διαχωρισμός συστάδων με euclidean για την πρώτη έως και την τέταρτη στήλη	9
2.10	Διαχωρισμός συστάδων με cityblock για την πρώτη έως και την τέταρτη στήλη	9
2.11	Διαχωρισμός συστάδων με cosine για την πρώτη έως και την τέταρτη στήλη	10
2.12	Διαχωρισμός συστάδων με correlation για την πρώτη έως και την τέταρτη στήλη	10
2.13	Διαχωρισμός συστάδων με euclidean για την πρώτη έως και την τρίτη στήλη	11
2.14	Διαχωρισμός συστάδων με cityblock για την πρώτη έως και την τρίτη στήλη	11
2.15	Διαχωρισμός συστάδων με cosine για την πρώτη έως και την τρίτη στήλη	12
2.16	Διαχωρισμός συστάδων με correlation για την πρώτη έως και την τρίτη στήλη	12
2.17	Διαχωρισμός συστάδων με euclidean για την δεύτερη έως και την τέταρτη στήλη	13
2.18	Διαχωρισμός συστάδων με cityblock για την δεύτερη έως και την τέταρτη στήλη	13
2.19	Διαχωρισμός συστάδων με cosine για την δεύτερη έως και την τέταρτη στήλη	14
2.20	Διαχωρισμός συστάδων με correlation για την δεύτερη έως και την τέταρτη στήλη	14
2.21	Διαχωρισμός συστάδων με βάση τις πρώτες δύο στήλες	16
2.22	Διαχωρισμός συστάδων με βάση όλες τις στήλες	16

2.23 Διαχωρισμός συστάδων με βάση τις στήλες 296 και 305	17
--	----

1 Σύνολο δεδομένων του Enron

Στην ακόλουθη ενότητα, θα γίνει συσταδοποίηση στο σύνολο δεδομένων των email του Enron.

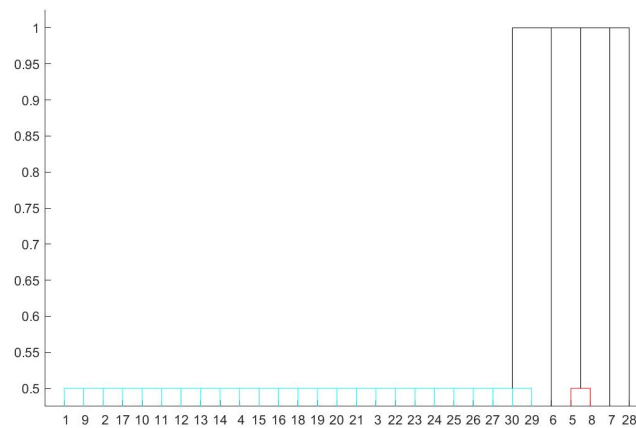
1.1 Δενδρόγραμμα για διαφορετικούς δεσμούς και δείκτες ομοιότητας

Το script 1.1, κάνει συσταδοποίηση της δεύτερης και τρίτης στήλης των δεδομένων με:

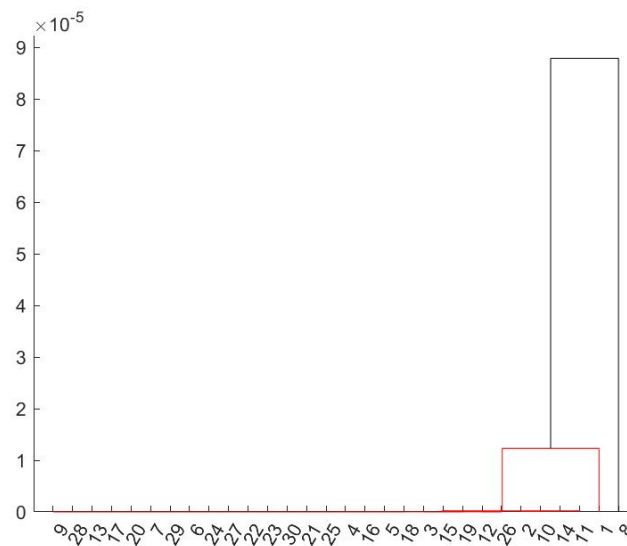
- απλό δεσμό και ομοιότητα jaccard
- απλό δεσμό και ομοιότητα cosine
- μέσο δεσμό και ομοιότητα cosine

και εμφανίζει το δενδρόγραμμα για κάθε ένα από αυτά.

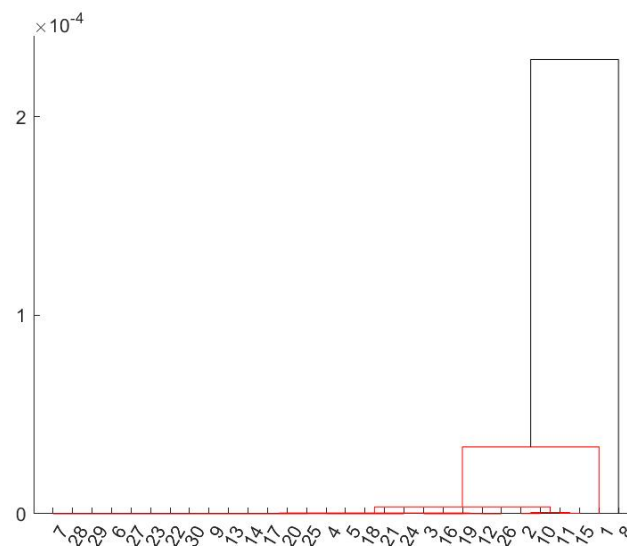
Είναι εμφανές σε ποίο επίπεδο θα γίνει η διαφοροποίηση σε κάθε ένα από τα τρία δενδρογράμματα. Ωστόσο, στο τρίτο δενδρόγραμμα (εικ. 1.3), αντί για τρεις συστάδες θα μπορούσαν να υπάρχουν τέσσερις. Επίσης, η μέθοδος που φαίνεται πως χωρίζει καλύτερα τα δεδομένα που έχουν περισσότερα κοινά, είναι με τον μέσο δεσμό και την ομοιότητα cosine.



Σχήμα 1.1: Απλός δεσμός με ομοιότητα jaccard



Σχήμα 1.2: Απλός δεσμός με ομοιότητα cosine



Σχήμα 1.3: Μέσος δεσμός με ομοιότητα cosine

```

1 load('enron100.mat')
2
3 d = pdist(en2(1:100,2:3),'jaccard');
4 Z = linkage(d); % single
5 [H,T] = dendrogram(Z,'ColorThreshold', 'default');
6

```

```
7 %print('images/single_jaccard','-djpeg')
8
9 figure
10 d = pdist(en2(1:100,2:3),'cosine');
11 Z = linkage(d); % single
12 [H,T] = dendrogram(Z,'ColorThreshold', 'default');
13 xtickangle(60);
14
15 %print('images/single_cosine','-djpeg')
16
17 figure
18 d = pdist(en2(1:100,2:3),'cosine');
19 Z = linkage(d, 'average');
20 [H,T] = dendrogram(Z,'ColorThreshold', 'default');
21 xtickangle(60);
22
23 %print('images/average_cosine','-djpeg')
```

Κώδικας 1.1: Script συσταδοποίησης των email του Enron

2 Διαμεριστική συσταδοποίηση με k-means

Σε αυτήν την ενότητα, θα γίνει συσταδοποίηση με k-means για διάφορες μεθόδους στα δεδομένα Iris και xV.

2.1 Διαχωριστική συσταδοποίηση για το σύνολο δεδομένων Iris

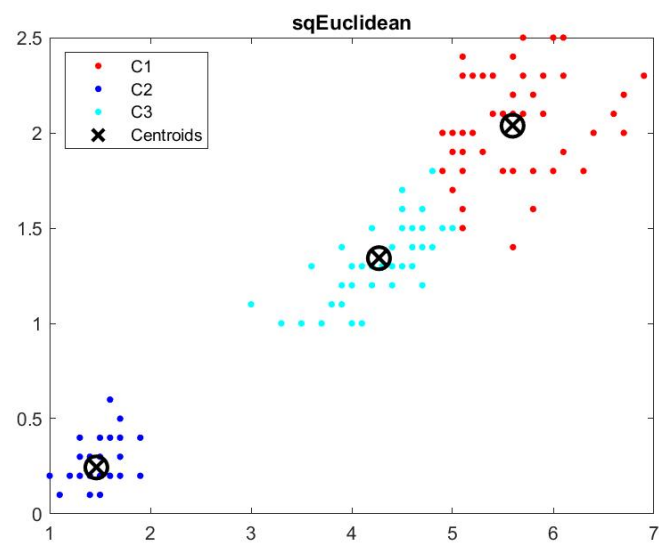
Στο script 2.1, υπολογίζεται ο k-means με τους εξής διαφορετικούς τρόπους υπολογισμού:

- Euclidean
- cityBlock
- cosine
- correlation

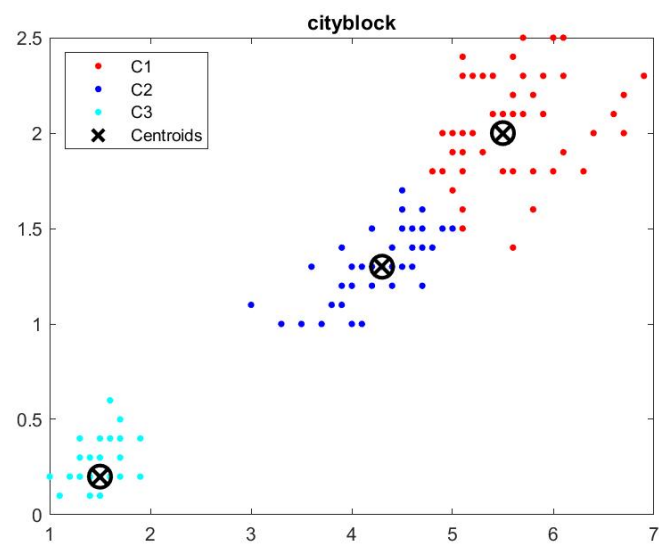
για κάθε ένα από τους εξής συνδιασμούς στηλών των δεδομένων Iris:

- τρίτη και τέταρτη
- πρώτη και δεύτερη
- πρώτη, δεύτερη, τρίτη και τέταρτη
- πρώτη, δεύτερη και τρίτη
- δεύτερη, τρίτη και τέταρτη

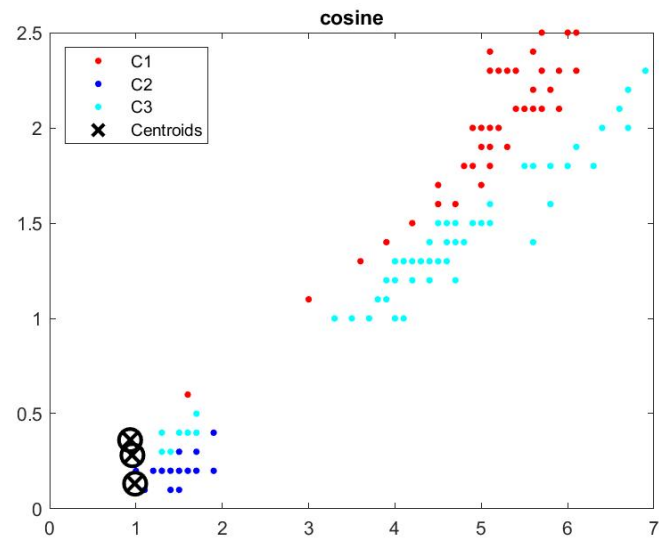
Από όλες τις διαφορετικές συσταδοποιήσεις, αυτή που φαίνεται να είναι η ποίο αποτελεσματική είναι η Euclidean και η CityBlock. Αυτό μπορούμε να το διακρίνουμε καλύτερα κάνοντας plot το IDX που είναι η πρώτη τιμή από τις δύο όπου επιστρέφει η kmeans. Ποίο συγκεκριμένα, εφόσον ξέρουμε ότι είναι χωρισμένα ανά 50, είναι φυσιολογικό να έχει γίνει η συσταδοποίηση αντίστοιχα στα πρώτα 50, στα επόμενα 50 και στα τελευταία 50 στοιχεία των δεδομένων.



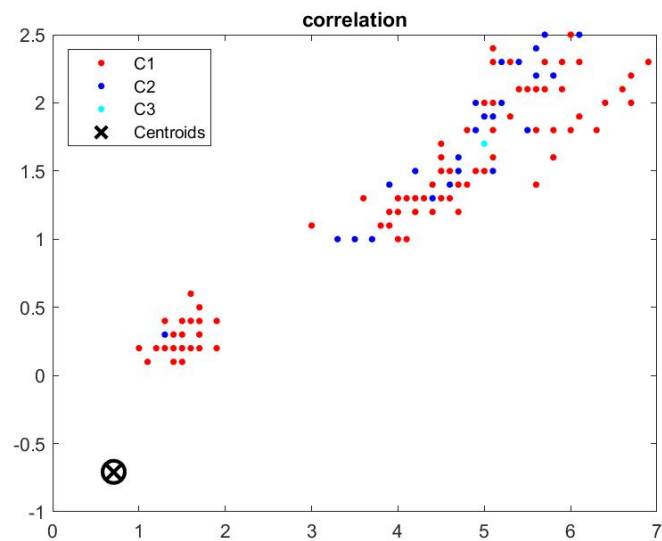
Σχήμα 2.1: Διαχωρισμός συστάδων με euclidean για την τρίτη και τέταρτη στήλη



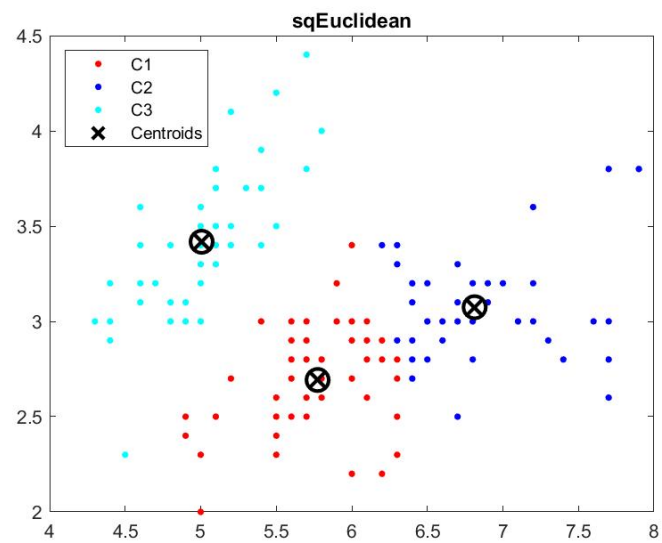
Σχήμα 2.2: Διαχωρισμός συστάδων με cityblock για την τρίτη και τέταρτη στήλη



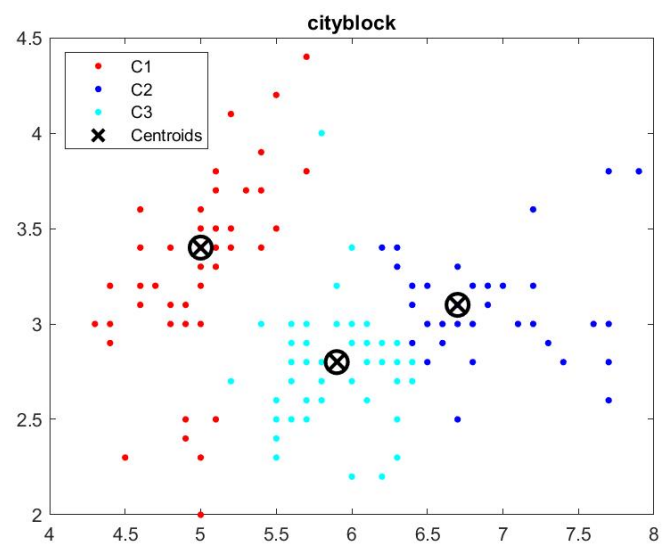
Σχήμα 2.3: Διαχωρισμός συστάδων με cosine για την τρίτη και τέταρτη στήλη



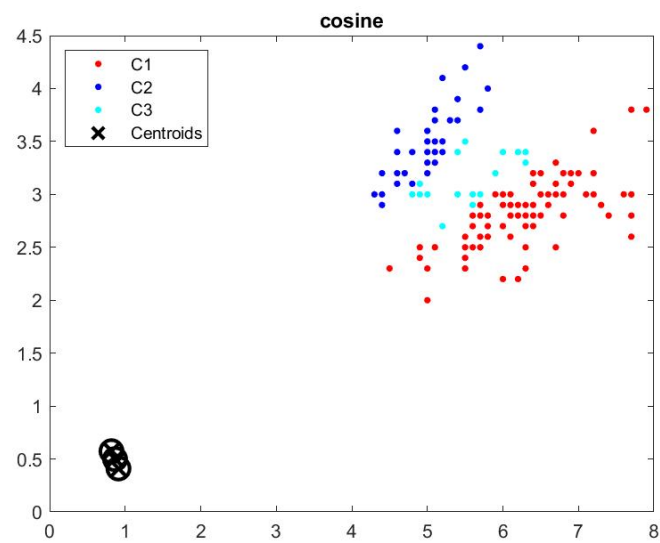
Σχήμα 2.4: Διαχωρισμός συστάδων με correlation για την τρίτη και τέταρτη στήλη



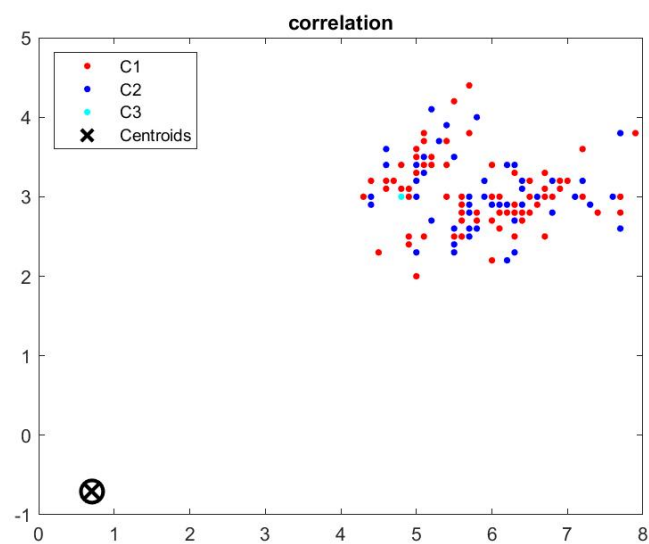
Σχήμα 2.5: Διαχωρισμός συστάδων με euclidean για την πρώτη και δεύτερη στήλη



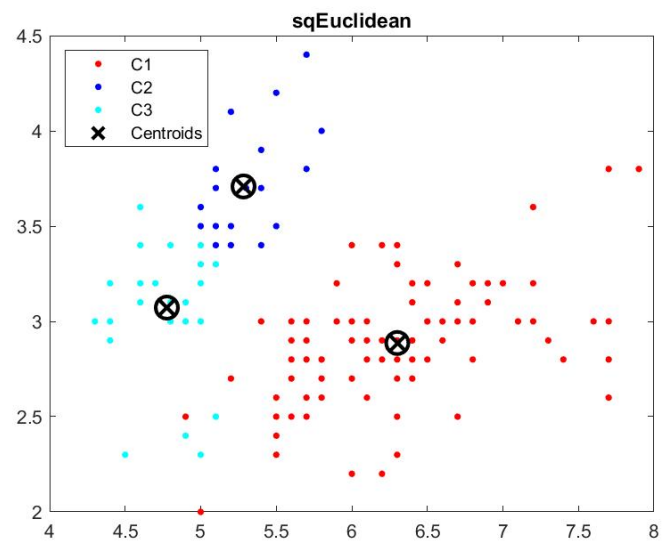
Σχήμα 2.6: Διαχωρισμός συστάδων με cityblock για την πρώτη και δεύτερη στήλη



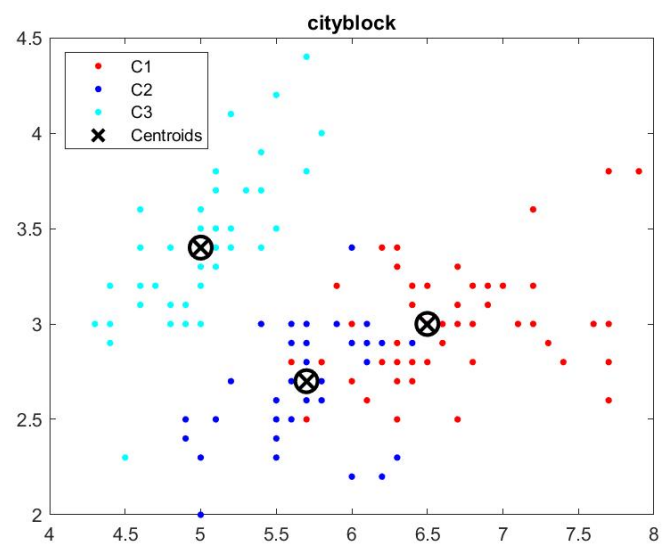
Σχήμα 2.7: Διαχωρισμός συστάδων με cosine για την πρώτη και δεύτερη στήλη



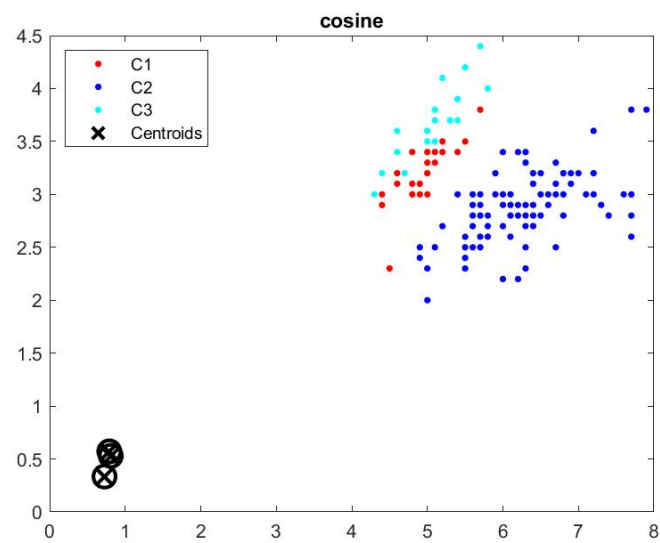
Σχήμα 2.8: Διαχωρισμός συστάδων με correlation για την πρώτη και δεύτερη στήλη



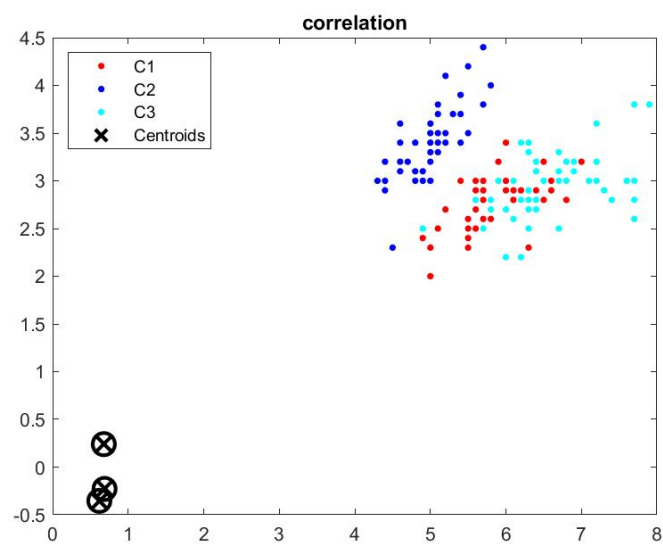
Σχήμα 2.9: Διαχωρισμός συστάδων με euclidean για την πρώτη έως και την τέταρτη στήλη



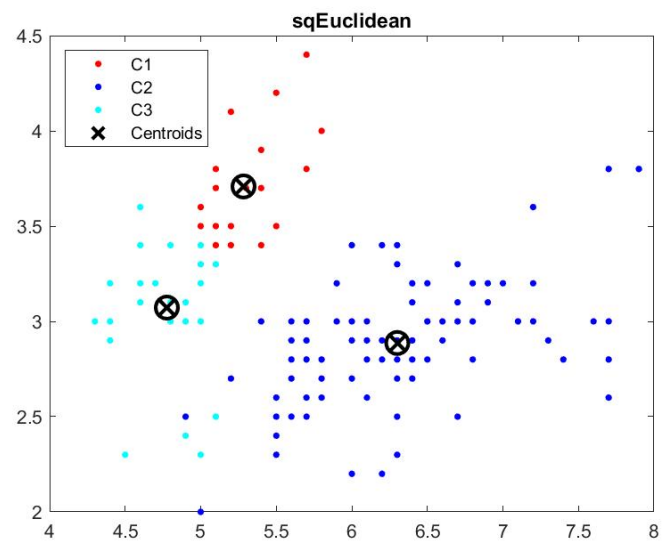
Σχήμα 2.10: Διαχωρισμός συστάδων με cityblock για την πρώτη έως και την τέταρτη στήλη



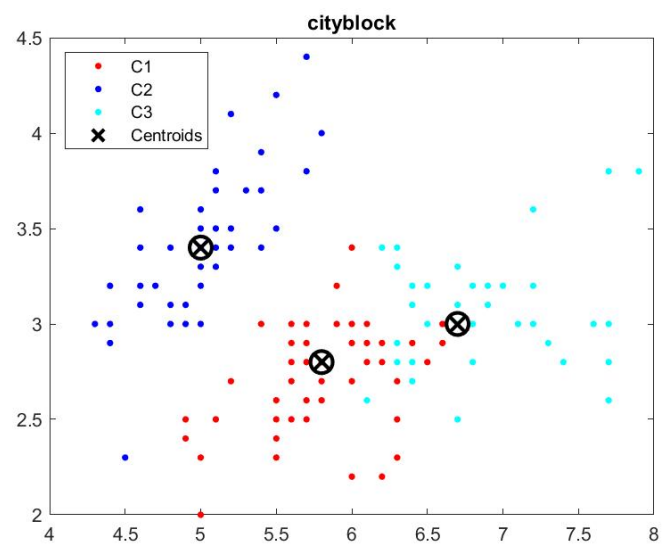
Σχήμα 2.11: Διαχωρισμός συστάδων με cosine για την πρώτη έως και την τέταρτη στήλη



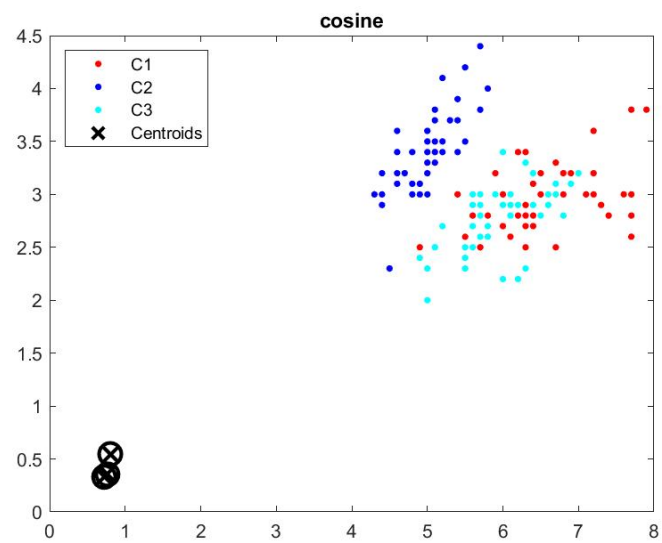
Σχήμα 2.12: Διαχωρισμός συστάδων με correlation για την πρώτη έως και την τέταρτη στήλη



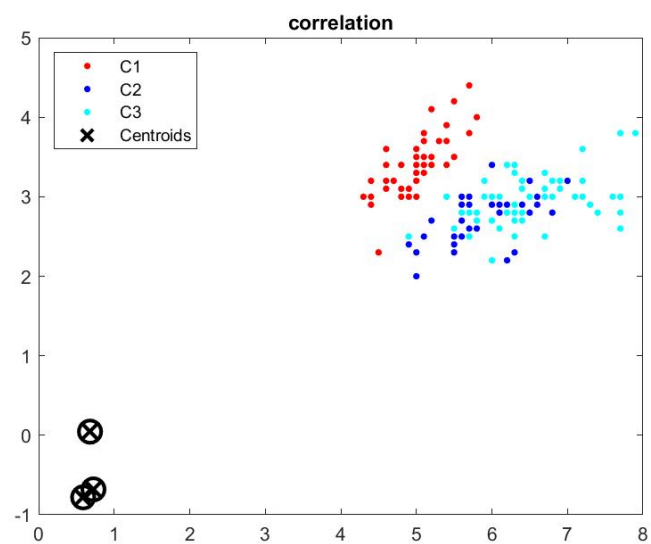
Σχήμα 2.13: Διαχωρισμός συστάδων με euclidean για την πρώτη έως και την τρίτη στήλη



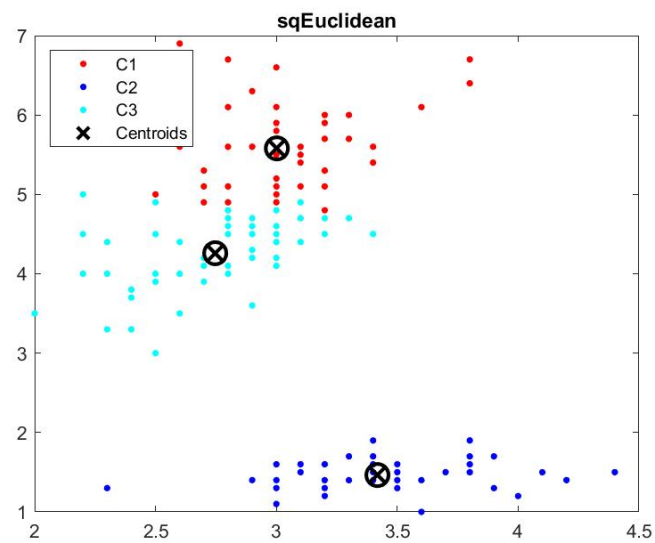
Σχήμα 2.14: Διαχωρισμός συστάδων με cityblock για την πρώτη έως και την τρίτη στήλη



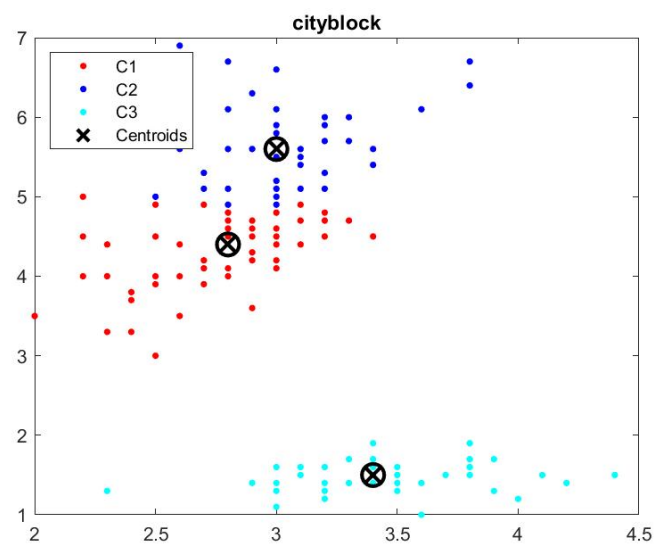
Σχήμα 2.15: Διαχωρισμός συστάδων με cosine για την πρώτη έως και την τρίτη στήλη



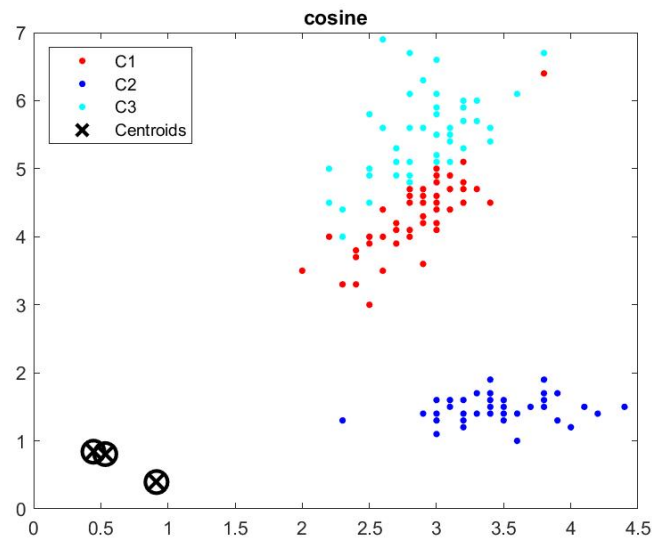
Σχήμα 2.16: Διαχωρισμός συστάδων με correlation για την πρώτη έως και την τρίτη στήλη



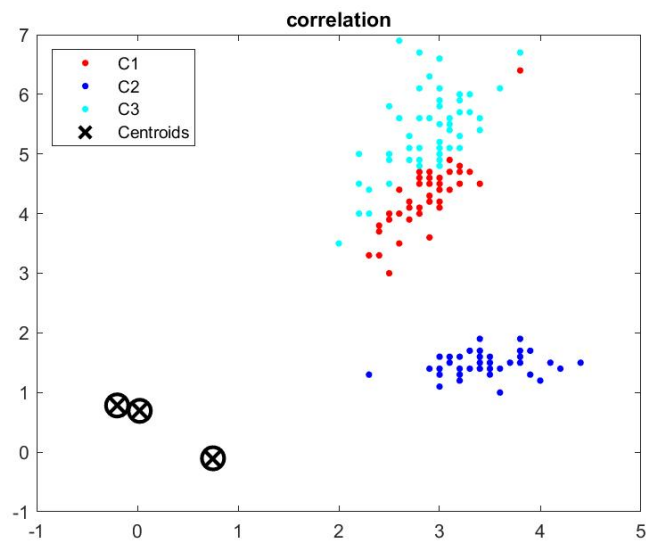
Σχήμα 2.17: Διαχωρισμός συστάδων με euclidean για την δεύτερη έως και την τέταρτη στήλη



Σχήμα 2.18: Διαχωρισμός συστάδων με cityblock για την δεύτερη έως και την τέταρτη στήλη



Σχήμα 2.19: Διαχωρισμός συστάδων με cosine για την δεύτερη έως και την τέταρτη στήλη



Σχήμα 2.20: Διαχωρισμός συστάδων με correlation για την δεύτερη έως και την τέταρτη στήλη

```

1 load('iris.dat')
2
3 method = {'sqEuclidean', 'cityblock',...
4           'cosine', 'correlation'};
5 plotColors = {'r', 'b', 'c'};
6 combinations = [3,4; 1,2; 1,4; 1,3; 2,4];

```

```

7 k = 3;
8
9 for j = 1:length(combinations)
10     X = iris(:,[combinations(j,1):combinations(j,2)]);
11
12     for i = 1:length(method)
13         [IDX, C] = kmeans(X,k, 'distance', method(i));
14
15         figure
16         plot(IDX,'o')
17
18         figure
19         for z=1:k
20             plot( X(IDX==z,1), X(IDX==z,2),...
21                 'LineStyle','none',...
22                 'Marker','.',...
23                 'color',plotColors{z},...
24                 'MarkerSize',12)
25             hold on
26         end
27
28         title(method(i));
29
30         plot(C(:,1),C(:,2),'kx', 'MarkerSize', 12, 'LineWidth', 2)
31         plot(C(:,1),C(:,2),'ko', 'MarkerSize', 12, 'LineWidth', 2)
32         legend('C1','C2','C3','Centroids','Location','NW')
33
34 %     filename = strcat('images/iris/comb(',...
35 %         num2str(combinations(j,1),'%d'), ',',...
36 %         num2str(combinations(j,2),'%d'), '),'_','...',
37 %         method{i});
38 %     print(filename,'-djpeg')
39 end
40 end

```

Κώδικας 2.1: Script συσταδοποίησης με k-means για διάφορους μεθόδους στα δεδομένα Iris

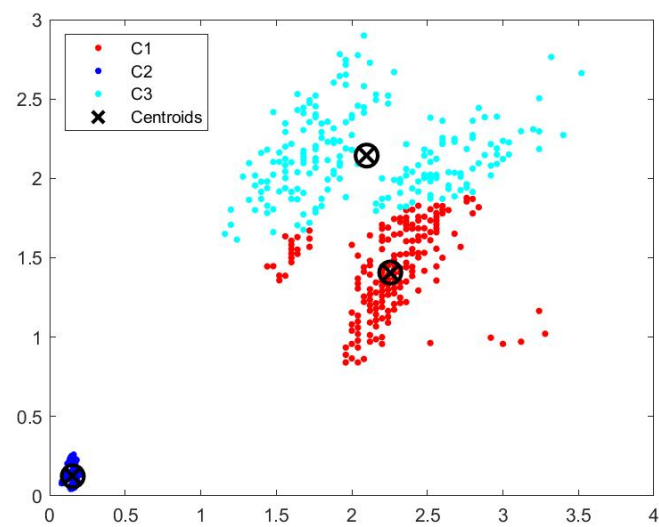
2.2 Διαχωριστική συσταδοποίηση για το σύνολο δεδομένων xV

Στην εικόνα 2.21, η συσταδοποίηση που έχει γίνει είναι αποδεκτή. Ωστόσο, οπτικά φαίνεται πως υπάρχει καλύτερος διαχωρισμός συσταδοποίησης για τα C1 και C3.

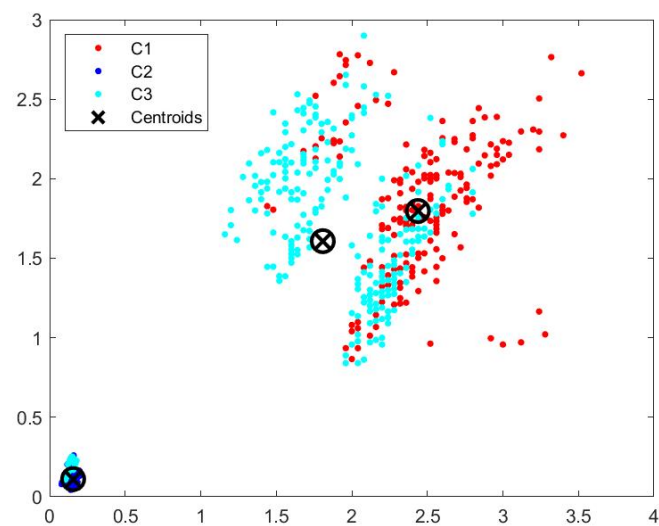
Ο διαχωρισμός με βάση όλες τις στήλες (εικ. 2.22) δεν είχε θεμιτό αποτέλεσμα και η συσταδοποίηση που υπάρχει δεν είναι ικανοποιητική.

Τέλος, η συσταδοποίηση με βάση τις στήλες 296 και 305 (εικ. 2.23), δεν είναι καθόλου εμφανές ο διαχωρισμός. Ενδεχομένως, να ήταν καλύτερος ο διαχωρισμός αν ο κεντροειδής του C3 ήταν δεξιότερα και του C1 ποίο ψηλά και δεξιά, δηλαδή κοντά στο σημείο (0.2, 0.2).

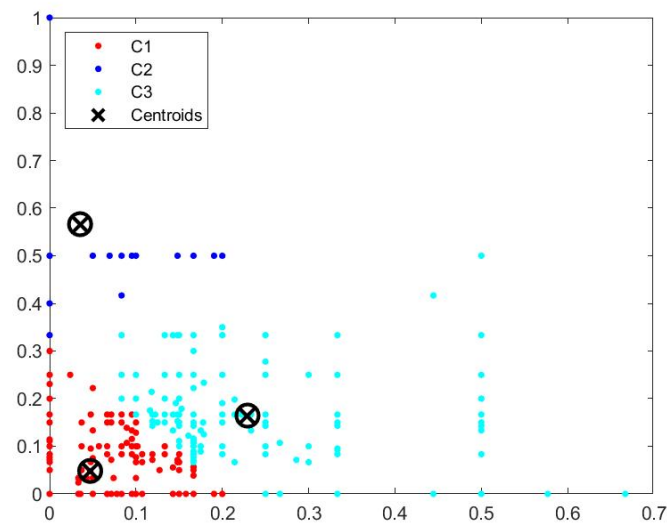
Είναι εμφανές, πως η καλύτερη συσταδοποίηση έγινε με τα δύο πρώτα χαρακτηριστικά (εικ. 2.21) όπου επέστρεψε οπτικά των καλύτερο διαχωρισμό.



Σχήμα 2.21: Διαχωρισμός συστάδων με βάση τις πρώτες δύο στήλες



Σχήμα 2.22: Διαχωρισμός συστάδων με βάση όλες τις στήλες



Σχήμα 2.23: Διαχωρισμός συστάδων με βάση τις στήλες 296 και 305

```

1 load('xV.mat')
2
3 notNaN = ~isnan(xV);
4 xV(~notNaN) = 0;
5
6 k = 3;
7 plotArgs = {'r', 'b', 'c'};
8
9 X = xV(:,1:2);
10 [IDX, C] = kmeans(X,k);
11
12
13 figure
14 for j=1:k
15     plot( X(IDX==j,1), X(IDX==j,2),...
16         'LineStyle','none',...
17         'Marker','.',...
18         'color',plotArgs{j},...
19         'MarkerSize',12)
20     hold on
21 end
22
23 plot(C(:,1),C(:,2),'kx', 'MarkerSize', 12, 'LineWidth', 2)
24 plot(C(:,1),C(:,2),'ko', 'MarkerSize', 12, 'LineWidth', 2)
25 legend('C1','C2','C3','Centroids','Location','NW')
26
27 % print('images/xV/comb(1,2)_sqEuclidean','-djpeg')
28

```

```
29 X = xV;  
30 [IDX, C] = kmeans(X,k);  
31  
32  
33 figure  
34 for j=1:k  
35     plot( X(IDX==j,1), X(IDX==j,2),...  
36         'LineStyle','none',...  
37         'Marker','.',...  
38         'color',plotArgs{j},...  
39         'MarkerSize',12)  
40     hold on  
41 end  
42  
43 plot(C(:,1),C(:,2),'kx', 'MarkerSize', 12, 'LineWidth', 2)  
44 plot(C(:,1),C(:,2),'ko', 'MarkerSize', 12, 'LineWidth', 2)  
45 legend('C1','C2','C3','Centroids','Location','NW')  
46  
47 % print('images/xV/comb(all)_sqEuclidean','-djpeg')  
48  
49  
50 X = [xV(:,296),xV(:,305)];  
51 [IDX, C] = kmeans(X,k);  
52  
53  
54 figure  
55 for j=1:k  
56     plot( X(IDX==j,1), X(IDX==j,2),...  
57         'LineStyle','none',...  
58         'Marker','.',...  
59         'color',plotArgs{j},...  
60         'MarkerSize',12)  
61     hold on  
62 end  
63  
64 plot(C(:,1),C(:,2),'kx', 'MarkerSize', 12, 'LineWidth', 2)  
65 plot(C(:,1),C(:,2),'ko', 'MarkerSize', 12, 'LineWidth', 2)  
66 legend('C1','C2','C3','Centroids','Location','NW')  
67  
68 % print('images/xV/comb(296,305)_sqEuclidean','-djpeg')
```

Κώδικας 2.2: Script συσταδοποίησης με k-means για διάφορους μεθόδους στα δεδομένα xV