

2

1

.

المستخلص

.

.

:

.

/

/

1

/

/

2

.....

:

Employment of Cluster Analysis and Nearest Neighbor Method in Pattern Recognition With the Application on Groundwater Quality in the Province of Nineveh

ABSTRACT

The main objective of pattern recognition methods is to develop capacity in the simulation of the most special objects in a different form, in order to classify these objects to a number of classes or types, these classes can be images or data or any kind of standards that need to be rated. In this research, recruitment of and the nearest neighbor method to the decision-making process in non-parametric pattern recognition through the use of theory in the Bayse theorem as a decision rule of discrimination, and their application in the field of groundwater quality in the province of Nineveh. Has given way suggest accurate results in the process of classifying data into two groups: the first water wells for drinking and other non-potable water, which could study the characteristics and qualities of each group and to identify qualities essential so it is easy to distinguish samples from any well in the province to see the quality of groundwater for The well under study.

Keyword: pattern recognition, nearest neighbor method, cluster analysis, groundwater.

: (1)

Patterns Recognition

Artificial

Machine Learning

. *Intelligence*

Classifiers

Intelligent Automated Systems

Pattern Space

.

)

(

)

(

:

(2)

.(

)

.....

Bayes Theory

:Cluster Analysis (3)

()

Clusters

()

)

(....

Unsupervised

Classification

()

()

[Theodoridis and ()

.Koutroumbas, 2006]

(1.3) الخطوات الأساسية لعملية العنقدة *Basic Steps of Clustering Process*:

في عملية العنقدة يوصف كل نمط بواسطة مجموعة من السمات *Set of Features* والتي من الممكن أن نعبر عنها بقيم حقيقية. وفي مرحلة التعليم يقدم كل نمط كمتجه يسمى بمتجه السمات الأساسية *Basic Features Vector*، ويمكن إدراج الخطوات الأساسية لعملية العنقدة كما يلي [Theodoridis and Koutroumbas, 2006]:

- اختيار السمات الأساسية ***Basic Features Selection***:

يجب أن يتم اختيار السمات الأساسية بشكل صحيح بحيث تمثل قدر كبير من المعلومات المتعلقة بالنمط قيد الدراسة، أي إعداد السمات الأساسية قبل استخدامها في المراحل اللاحقة.

- معيار التعنقد ***Clustering Criterion***:

يعتمد معيار العنقدة على الخبرة المكتسبة للخبير مستند على نوع العناقيد التي يتوقع أنها تقع تحت مجموعة من البيانات، فقد يتضمن احد العناقيد

.....

مجموعة من متجهات ذات السمات الملائمة لمعيار معين، في حين أن عنقود آخر
يعتمد على معيار يختلف عن المعيار الأول في تعنقه.

:Clustering Algorithms -

:Validity of Results -

:Interpretation of the Results -

: K-Nearest Neighbour Classification (4)

(closest)

.continuous

categorical

Probability Theory

Naïve Bayes Algorithm

.[Bramer,2007]

()

k_i

k

k

.[Bramer,2007]

[Webb,2002]:

(1.4)

x

: v

x

v

$$\theta = \int_{v(x)} p(x) dx \quad \dots(1)$$

:

v

$$\theta \sim p(x)v \quad \dots(2)$$

v

θ

v

n

()

k

x

:

$$p(x|\omega_m) \cdot p(\omega_m) > p(x|\omega_i) \cdot p(\omega_i) \quad \dots(7)$$

(4)

$p(x)$

ω_m k_m k

$$\sum_{m=1}^c k_m = k \quad \dots(8)$$

n_m ω_m

$$\sum_{m=1}^c n_m = n \quad \dots(9)$$

Class-Conditional Density

$$\hat{p}(x|\omega_m) = \frac{k_m}{n_m \cdot v} \quad \dots(10)$$

Prior Probability

$$\hat{p}(\omega_m) = \frac{n_m}{n} \quad \dots(11)$$

(7) (11) (10)

$$\frac{k_m}{n_m \cdot v} \cdot \frac{n_m}{n} > \frac{k_i}{n_i \cdot v} \cdot \frac{n_i}{n} \quad \forall i \in c \quad \dots(12)$$

ω_m (class) x

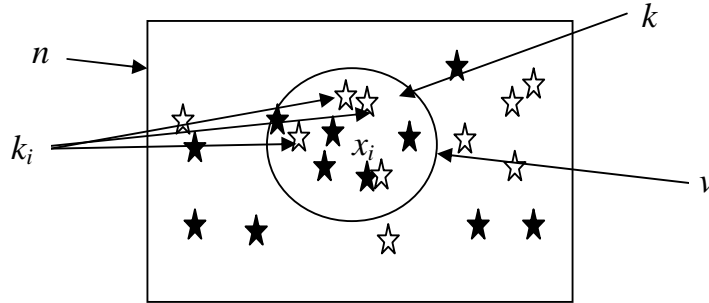
$$k_m > k_i \quad \forall i \in c \quad \dots(13)$$

Discriminate Function

v

:

$$h(x) = \frac{k_i}{k} \quad \forall i \in c \quad \dots(14)$$



الشكل رقم (1) : يوضح خوارزمية الجار الأقرب

:Similarity and Distance Measures

(5)

[Bramer,2007]

() ()

Dissimilarity

$$dist(x_i, x_j)$$

Distance

n

Euclidean Distance

:

$$dist(x_i, x_j) = \sqrt{\sum_{t=1}^n (x_{it} - x_{jt})^2} \quad \dots(15)$$

$$. \text{dist}(x_i, x_i) = 0 \quad .1$$

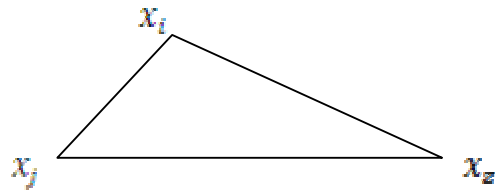
$$: \quad x_i \quad x_j \quad \quad \quad x_j \quad x_i \quad .2$$

$$\text{dist}(x_i, x_j) = \text{dist}(x_j, x_i)$$

$$: \quad .3$$

$$\text{dist}(x_i, x_j) \leq \text{dist}(x_i, x_2) + \text{dist}(x_2, x_j)$$

:



الشكل رقم (2) : يوضح متباينة المثلث

: (6)

)

(

(Recharge Zones)

.....

)

(Fresh

[1989]

)

(

[Mahmood, 1994]

()

[2001]

[2004] .

: (1.6)

1998

[*APHA, AWWA, WEF, 1998*]

: ()

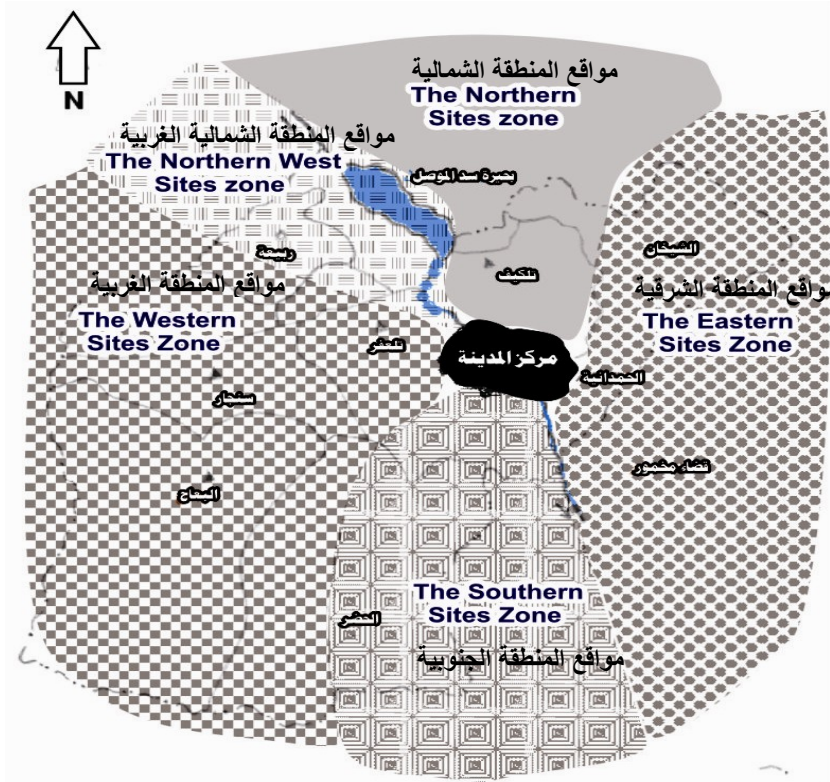
.() -1

. -2

. -3

. -4

:(3)



: (3)

()

:(1)

code	نوع التركيز	العسرة الكلية	المواد الصلبة الكلية	المواد الصلبة العالقة	المواد الصلبة الذاتية	الدالة الحامضية	التوصيلة الكهربائية	عسرة الكالسيوم	الكلوريدات	القاعدية	الفوسفات	النترات	الكبريتات
	المنطقة	T.H	T.S	T.D.S	S.S	PH	EC	Ca.H	Cl ⁻	Alk.	PO ₄	NO ₃	SO ₄
1	شرق ربيعة	184	1129	699	430	8.45	1450	32	35	12	0.03	5.4	190
2	أم الكهف	1766	10164	4331	5833	7.82	7613	633	1311	7.3	0.04	7.8	476.6
3	قرية خرمل	600	1834	1560	274	7.72	2340	205	151.7	8	0.04	8.2	500
4	قرية عوينات	376	1555	1292	263	7.99	1780	160	167.7	6	0	14	560
5	العياضية	1350	4469	4251	218	7.97	7770	420	938	12	0.25	22	2150
6	بجاج خويتلة	2350	4524	4252	272	6.98	5500	1450	275.5	10	0.29	27	2520
7	ربيعة قرية السعدة	390	682	650	32	8.6	1850	130	147.6	10	0	8.6	210
8	وادي عكاب موصل	2000	3274	3000	274	6.5	3560	1500	179.4	80	0.38	0.5	2138
9	زمار	500	873	843	30	7.7	1100	230	127.5	328	0.86	13.7	212
10	الخازر	481.5	754	742	12	7.3	1079	152	571	14	0.08	20.4	141
11	الشاقولي	561.5	1293	1257.5	35.5	7.37	2850	405.5	159	240	0.046	0	446
12	الشلالات	1612	2024	1952	72	7.13	2330	1224	13	12	0.16	0	750
13	بعويزة	2840	5708	5196	512	6.85	7300	710	441	30	0	14.4	3300
14	برطلة	1460	3131	3120	11	6.8	3470	908	370	110	0.05	14.5	1700
15	حقل دواجن برطلة	2385	4583	4384	200	6.98	6560	1085	642	275	0.02	14.3	2037
16	مخمور	480	830	754	76	6.88	1025	288	31	152	0.03	0	375
17	كوكجلي	1150	1863	1709	154	6.89	3690	570	238.5	468	0.04	6.2	692
18	النمرود	2240	10996	5350	5646	7.3	7450	925	511.7	287.5	0.15	13.89	2589
19	مفرق حمام العليل	2440	3082	3021	61	7.4	3250	1730	69	255	2.85	0.16	1850
20	القوش	286	579	500	79	7.66	538	162	4.7	300	0.05	2.25	36
21	فايدة	270	396	390	6	7.98	712	116	17	284	0.2	6.75	17

.....

: (2.6)

Minitab(14)

(variable)

.(observation)

(Ward's Method)

(Agglomerative Method)

: *(two classes) (cluster observation)*

Cluster Analysis of Observations:

T.H; T.S ; T.D.S ; S.S ; PH ; EC ; Ca.H ; Cl⁻ ; Alk. ; PO₄ ; NO₃ ; SO₄

Standardized Variables, Euclidean Distance, Ward Linkage

Amalgamation Steps

Step	Number of clusters	Similarity level	Distance level	Clusters joined		New cluster	Number of obs. in new cluster
1	20	90.932	0.7379	1	7	1	2
2	19	89.426	0.8604	20	21	20	2
3	18	88.119	0.9667	3	4	3	2
4	17	82.754	1.4033	11	16	11	2
5	16	73.532	2.1537	9	20	9	3
6	15	72.638	2.2264	8	12	8	2
7	14	71.172	2.3456	1	3	1	4
8	13	70.670	2.3865	14	15	14	2

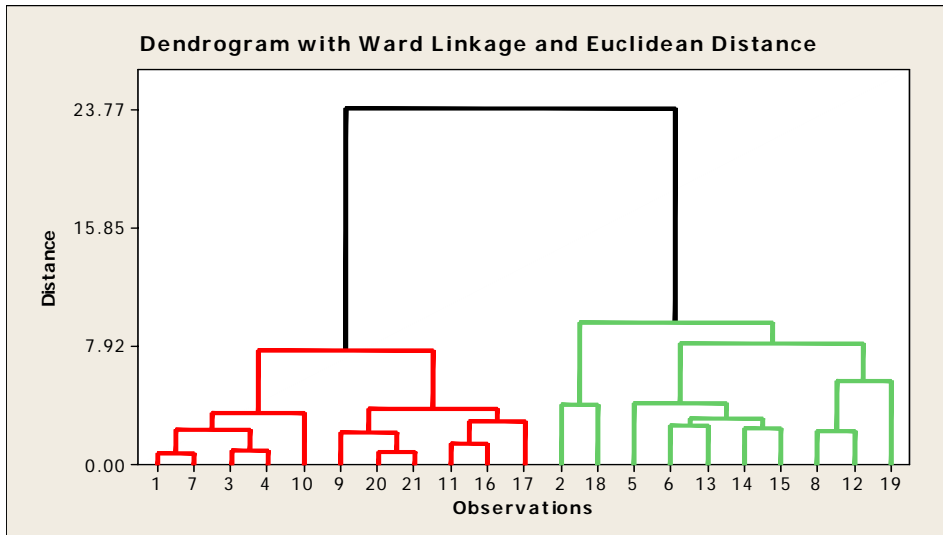
9	12	67.526	2.6424	6	13	6	2
10	11	64.437	2.8937	11	17	11	3
11	10	61.970	3.0944	6	14	6	4
12	9	57.813	3.4326	1	10	1	5
13	8	54.503	3.7020	9	11	9	6
14	7	50.602	4.0194	2	18	2	2
15	6	50.236	4.0492	5	6	5	5
16	5	30.993	5.6149	8	19	8	3
17	4	6.303	7.6239	1	9	1	11
18	3	0.017	8.1354	5	8	5	8
19	2	-16.883	9.5105	2	5	2	10
20	1	-192.168	23.7730	1	2	1	21

Final Partition

Number of clusters: 2

	Number of observations	Within cluster sum of squares	Average distance from centroid	Maximum distance from centroid
Cluster1	11	37.761	1.78676	2.91990
Cluster2	10	106.017	3.08735	4.88369

: (Dendrogram)



(Dendrogram) : (4)

()

:

1. (ω_1)
 . (7.6239)
2. (ω_2)
 . (9.5105)

:

:(2)

العنقود الأول ω_1	العنقود الثاني ω_2	
شرق ربيعة	أم الكهف	.1
قرية خرمل	العياضية	.2
قرية عوينات	بعاج خويتلة	.3
ربيعة قرية السعدة	وادي عكاب موصل	.4
زمار	الشلالات	.5
الخازر	بعويزة	.6
الشاقولي	برطلة	.7
مخمور	حقل دواجن برطلة	.8
كوكجلي	النمرود	.9
القوش	مفرق حمام العليل	.10
فايدة		.11

(SO₄)

(/ 250)

.[WHO,1985]

(/ 500)

.

() (T.H)

500)

(/

.[WHO,1985]

:

_____ (ω_1) .1

.
_____ (ω_2) .2

\underline{x}_0 : (3)

المنطقة	المسافة Distance	رمز المنطقة Code	التصنيف للمناطق Classification
بعاج خويتلة	1483.5	1	ω_2
وادي عكاب موصل	1529.5	2	ω_2
برظلة	1848.4	3	ω_2
مفرق حمام العليل	1933.2	4	ω_2
حقل دواجن برظلة	2213.5	5	ω_2
كوكجلي	3342.6	6	ω_1
العياضية	3374.9	7	ω_2
الشلالات	3680.2	8	ω_2
بعويزة	4012	9	ω_2
قرية خرمل	4425.8	10	ω_1
الشاقولي	4518.1	11	ω_1
قرية عوينات	5062	12	ω_1
ربيعة قرية السعدة	5815.9	13	ω_1
شرق ربيعة	5914.8	14	ω_1
زمار	6036.8	15	ω_1
مخمور	6086.5	16	ω_1
الخازر	6190.5	17	ω_1
القوش	6801.4	18	ω_1
فايدة	6840.7	19	ω_1
أم الكهف	9347.8	20	ω_2
النمرود	9775	21	ω_2

9-Nearest Neighbour Classification

(8)

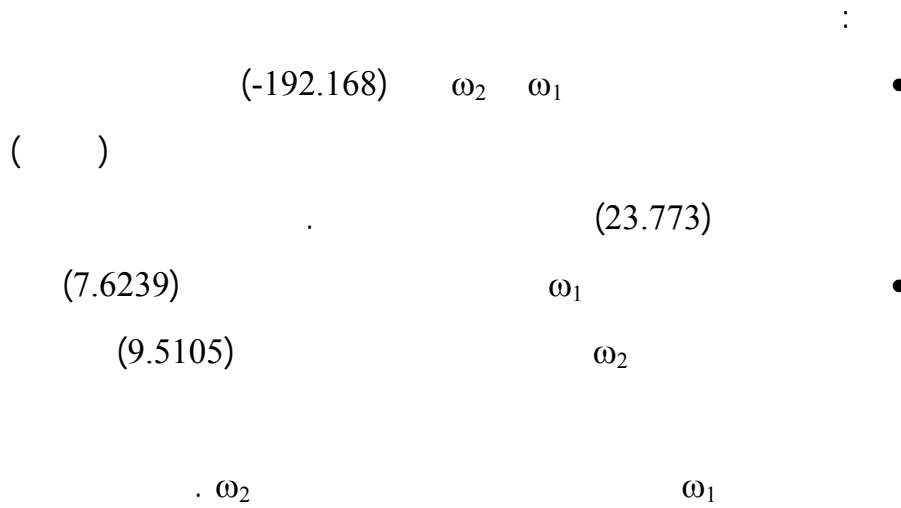
(k=9)

$$\begin{aligned}
 & \dots\dots\dots \\
 & \quad \quad \quad (\quad \quad \quad) \omega_2 \\
 & \quad \quad \quad : \quad \quad \quad (\quad \quad \quad) \omega_1 \\
 & h_{\omega_1}(\underline{x}_0) = \frac{1}{9} = 11.11\% \quad , \quad h_{\omega_2}(\underline{x}_0) = \frac{8}{9} = 88.89\% \\
 & \quad \quad \quad 88.89\% \\
 & \quad \quad \quad 11.11\% \quad \quad \quad \underline{x}_0
 \end{aligned}$$

: (7)

(Ward's Method) .1
 (cluster observation) (Agglomerative Method)

(4)



.2

.3

:

.1

:

.2

المصادر

" (2006) .1

"

" (2004) .2

"

" (2001) .3

.....

" (1989)

.4

. 20-11 9

"

5. APHA,AWWA,WPCF (1985), "*Standard Methods for Water Examination and Tests*" New-York .
6. Bramer, M. (2007), "*Principles of Data Mining*", Springer-Verlag London Limited.
7. Mahmood, F.Y. (1994), "*Physio-chemical evaluation of groundwater of some wells in Nineveh district used for drinking and domestic purpose*" . Al-Rafidain engineering journal. Vol. 2, No.2 .
8. Theodoridis S. and Koutroumbas K. (2006), "*Pattern Recognition*", Third edition, Elsevier, USA.
9. Webb, Andrew R.(2002), " *Statistical Pattern Recognition*", Second Edition, John Wiley & Sons, Ltd.
10. WHO (1985), "*Guidelines for drinking water quality* " Vol. 3, Geneva.

الملحق

برنامج لقياس المسافة بين متجه السمات x_0 الخاص بعينة أحد الآبار في المحافظة
وبين جميع الآبار في المحافظة ضمن الفنتين قيد الدراسة

```
clear
clc
x1=[184.00 1766.0 600.00 376.00 1350.00 2350.00 390.0
2000.00 500.00 481.50 561.50 1612.00 2840.00 1460.00 2385.00
480.00 1150.00 2240.0 2440.00 286.00 270.00];
x2=[1129.00 10164.0 1834.00 1555.00 4469.00 4524.00 682.0
3274.00 873.00 754.00 1293.00 2024.00 5708.00 3131.00 4583.00
830.00 1863.00 10996.0 3082.00 579.00 396.00];
x3=[699.00 4331.0 1560.00 1292.00 4251.00 4252.00 650.0
3000.00 843.00 742.00 1257.50 1952.00 5196.00 3120.00 4384.00
754.00 1709.00 5350.0 3021.00 500.00 390.00];
x4=[430.00 5833.0 274.00 263.00 218.00 272.00 32.0 274.00
30.00 12.00 35.50 72.00 512.00 11.00 200.00 76.00
154.00 5646.0 61.00 79.00 6.00];
x5=[8.45 7.8 7.72 7.99 7.97 6.98 8.6 6.50 7.70
7.30 7.37 7.13 6.85 6.80 6.98 6.88 6.89
7.3 7.40 7.66 7.98];
x6=[1450.00 7613.0 2340.00 1780.00 7770.00 5500.00 1850.0
3560.00 1100.00 1079.00 2850.00 2330.00 7300.00 3470.00 6560.00
1025.00 3690.00 7450.0 3250.00 538.00 712.00];
x7=[32.00 633.0 205.00 160.00 420.00 1450.00 130.0
1500.00 230.00 152.00 405.50 1224.00 710.00 908.00 1085.00
288.00 570.00 925.0 1730.00 162.00 116.00];
x8=[35.00 1311.0 151.70 167.70 938.00 275.50 147.6 179.40
127.50 571.00 159.00 13.00 441.00 370.00 642.00 31.00
238.50 511.7 69.00 4.70 17.00];
x9=[12.00 7.3 8.00 6.00 12.00 10.00 10.0 80.00
328.00 14.00 240.00 12.00 30.00 110.00 275.00 152.00
468.00 287.5 255.00 300.00 284.00];
x10=[0.03 0.0 0.04 0.00 0.25 0.29 0.0 0.38 0.86
0.08 0.05 0.16 0.00 0.05 0.02 0.03 0.04
0.2 2.85 0.05 0.20];
x11=[5.40 7.8 8.20 14.00 22.00 27.00 8.6 0.50 13.70
20.40 0.00 0.00 14.40 14.50 14.30 0.00 6.20
13.9 0.16 2.25 6.75];
x12=[190.00 476.6 500.00 560.00 2150.00 2520.00 210.0
2138.00 212.00 141.00 446.00 750.00 3300.00 1700.00 2037.00
375.00 692.00 2589.0 1850.00 36.00 17.00];

f=[2050 3590 3426 164 7.29 4980 1500 252 30 0.06 14.2 2000];

for i=1:21
    w=(f(1)-x1(i))^2+(f(2)-x2(i))^2+(f(3)-x3(i))^2+(f(4)-
x4(i))^2+(f(5)-x5(i))^2+(f(6)-x6(i))^2+(f(7)-x7(i))^2+(f(8)-
x8(i))^2+(f(9)-x9(i))^2+(f(10)-x10(i))^2+(f(11)-x11(i))^2+(f(12)-
x12(i))^2;
    d=sqrt(w)
end
```