

НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ
«КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ»
ФАКУЛЬТЕТ ІНФОРМАТИКИ ТА ОБЧИСЛЮВАЛЬНОЇ ТЕХНІКИ
Кафедра автоматизованих систем обробки інформації та управління

ЗВІТ

до лабораторної роботи №1
з дисципліни «Інтелектуальний аналіз даних»

Виконав:

студент гр. ІС-32
Капорін Р. М.

Мета: спрогнозувати погоду в Житомирі за 2013 рік.

	ЯНВ	ФЕВ	МАР	АПР	МАЙ	ІЮН	ІЮЛ	АВГ	СЕН	ОКТ	НОЯ	ДЕК
1	5	2	1	2	15	19	20	25	23	8	15	3
2	3	1	0	4	20	19	24	26	19	6	15	3
3	2	1	1	3	28	17	25	23	14	6	18	1
4	2	-1	-2	1	18	22	27	28	19	9	13	2
5	1	4	4	4	22	25	29	30	19	10	16	2
6	0	4	9	2	23	27	30	30	19	15	14	1
7	-5	3	5	2	26	26	26	31	21	16	9	0
8	-6	-1	1	6	27	27	24	32	20	16	16	-2
9	-8	-1	-3	8	28	27	28	33	23	13	14	2
10	-3	-2	-1	12	27	28	27	32	18	11	12	-4
11	-1	0	-2	12	29	24	30	21	22	16	7	1
12	-3	-2	0	13	28	20	22	27	21	13	6	1
13	-3	1	3	10	24	24	21	31	16	15	5	1
14	-7	1	1	11	17	26	22	22	14	13	5	0
15	1	-1	-1	10	22	26	21	23	14	14	5	0
16	1	-2	-6	17	25	27	22	24	18	13	4	1
17	0	-3	-4	18	27	26	23	26	22	14	5	4
18	-4	-4	-1	20	27	26	27	28	16	14	7	1
19	-5	-2	1	24	25	24	28	31	12	8	7	-1
20	-3	-1	2	20	29	27	20	30	11	17	12	0
21	-2	-1	3	14	23	29	20	31	14	18	11	2
22	-4	0	0	17	21	30	21	20	16	15	11	5
23	-7	1	-8	22	21	33	19	21	17	16	9	7
24	-6	-1	-4	19	18	26	21	21	14	21	9	7
25	-8	2	-4	22	19	29	24	19	13	14	9	2
26	-10	2	-2	26	21	24	24	17	7	15	-1	4
27	-12	4	-1	30	15	20	27	18	10	22	-4	4
28	-8	3	0	25	19	25	31	21	11	18	1	4
29	-1		2	24	23	21	32	18	11	22	2	11
30	1		6	24	25	23	32	18	12	13	4	4
31	3		2		20		24	23		12		1

Рис. 1 - Початкові дані

	1 january	2 february	3 march	4 april	5 may	6 june	7 july	8 august	9 september	10 october	11 november	12 december
1	5	2	1	2	15	19	20	25	23	8	15	3
2	3	1	0	4	20	19	24	26	19	6	15	3
3	2	1	1	3	28	17	25	23	14	6	18	1
4	2	-1	-2	1	18	22	27	28	19	9	13	2
5	1	4	4	4	22	25	29	30	19	10	16	2
6	0	4	9	2	23	27	30	30	19	15	14	1
7	-5	3	5	2	26	26	26	31	21	16	9	0
8	-6	-1	1	6	27	27	24	32	20	16	16	-2
9	-8	-1	-3	8	28	27	28	33	23	13	14	2
10	-3	-2	-1	12	27	28	27	32	18	11	12	-4
11	-1	0	-2	12	29	24	30	21	22	16	7	1
12	-3	-2	0	13	28	20	22	27	21	13	6	1
13	-3	1	3	10	24	24	21	31	16	15	5	1
14	-7	1	1	11	17	26	22	22	14	13	5	0
15	1	-1	-1	10	22	26	21	23	14	14	5	0
16	1	-2	-6	17	25	27	22	24	18	13	4	1
17	0	-3	-4	18	27	26	23	26	22	14	5	4
18	-4	-4	-1	20	27	26	27	28	16	14	7	1
19	-5	-2	1	24	25	24	28	31	12	8	7	-1
20	-3	-1	2	20	29	27	20	30	11	17	12	0
21	-2	-1	3	14	23	29	20	31	14	18	11	2
22	-4	0	0	17	21	30	21	20	16	15	11	5
23	-7	1	-8	22	21	33	19	21	17	16	9	7
24	-6	-1	-4	19	18	26	21	21	14	21	9	7
25	-8	2	-4	22	19	29	24	19	13	14	9	2
26	-10	2	-2	26	21	24	24	17	7	15	-1	4
27	-12	4	-1	30	15	20	27	18	10	22	-4	4
28	-8	3	0	25	19	25	31	21	11	18	1	4
29	-1		2	24	23	21	32	18	11	22	2	11
30	1		6	24	25	23	32	18	12	13	4	4
31	3		2		20		24	23		12		1

Рис. 2 - Початкові дані у Statistica

Тут змінним відповідають місяці року, а спостереженням - дні місяців. По своїй природі змінні неперервні, хоч на сайті "Ну і погода" й були представлені лише цілі значення температур. Метою нашого дослідження буде визначення можливих залежностей між змінними, побудова моделі для прогнозування температур по наявним даним.

ОПИСОВИЙ АНАЛІЗ

Побудуємо графік температур квітня за 2013 рік:

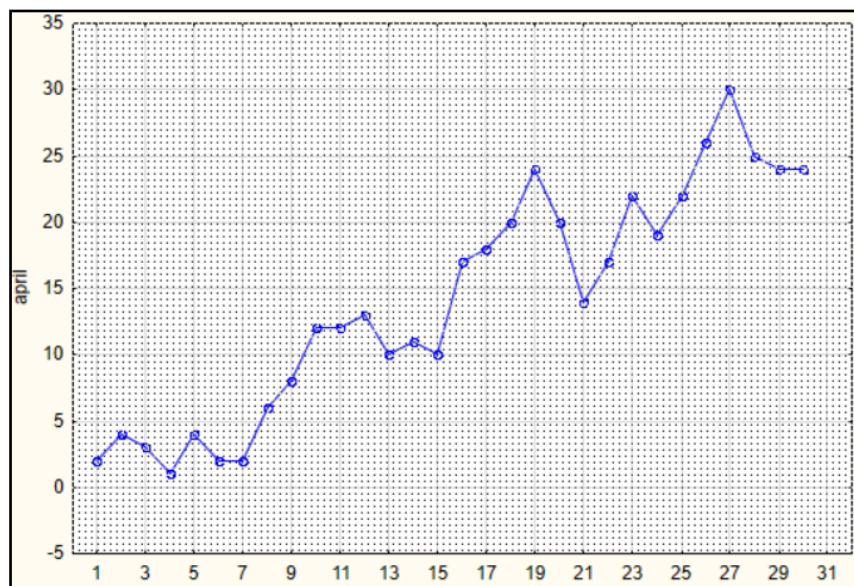


Рис. 3 - Графік температур квітня 2013

Бачимо, що температура протягом січня нерівно, але постійно зростає. Виконаємо те ж саме для січня:

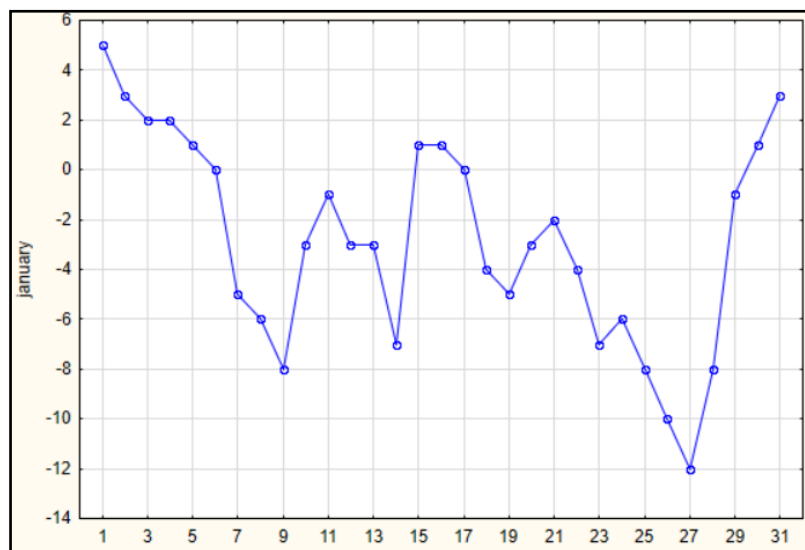


Рис. 4 - Графік температур січня 2013

Тепер побудуємо графіки всіх місяців кожної пори року. Результати представлені на рис. 5-8.

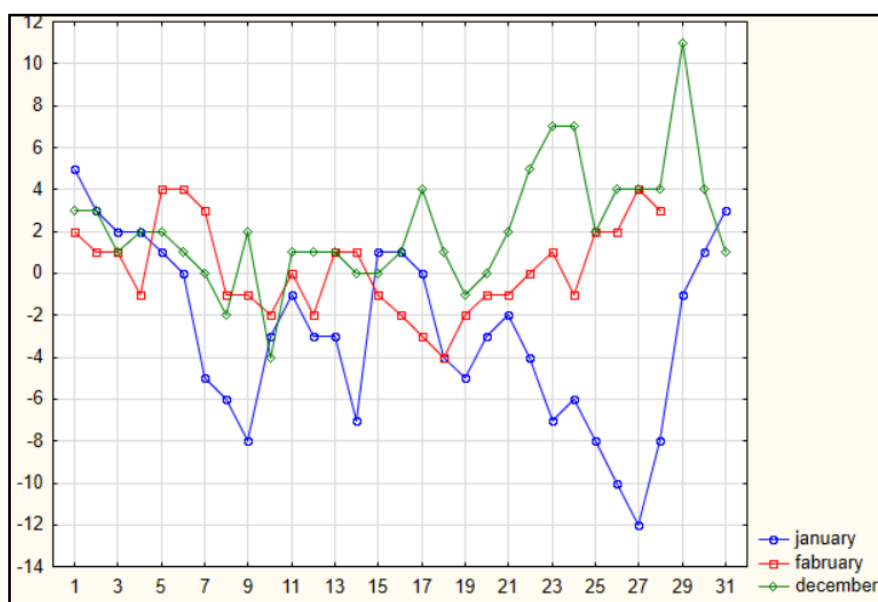


Рис. 5 - Графік температур зимових місяців

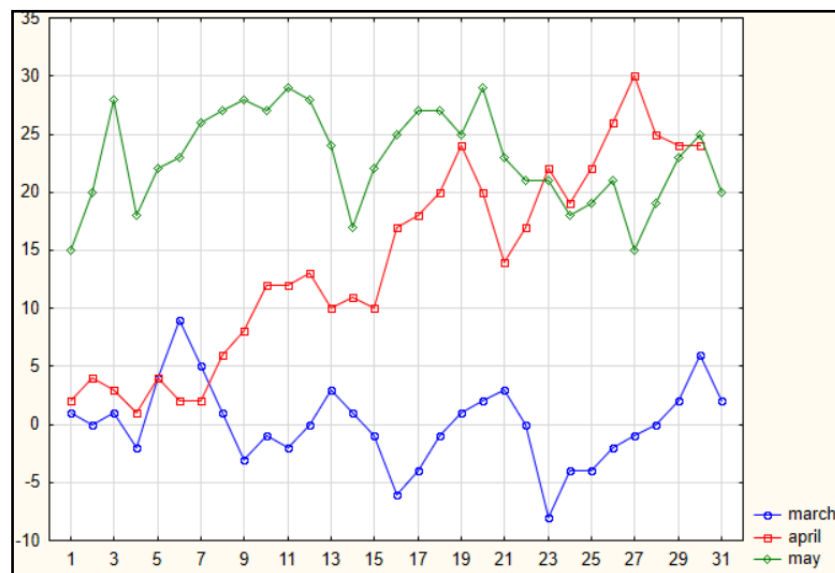


Рис. 6 - Графік температур весняних місяців

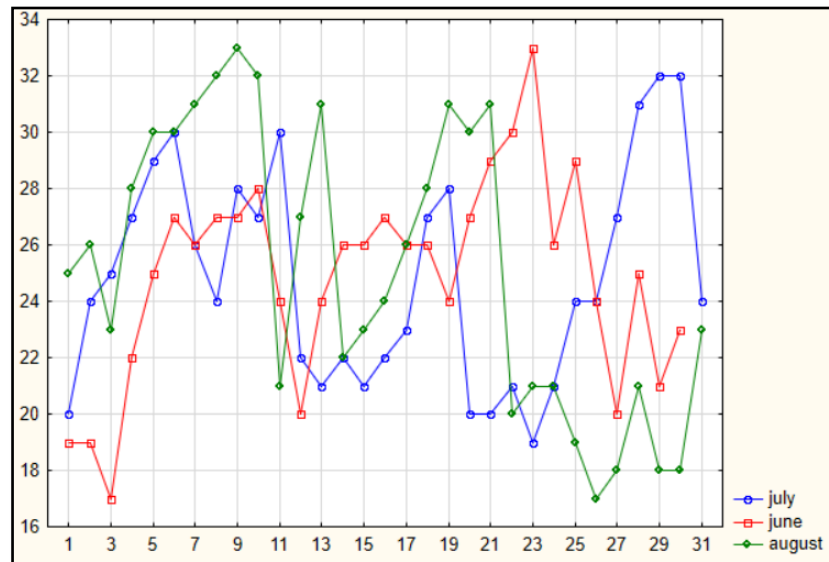


Рис. 7 - Графік температур літніх місяців

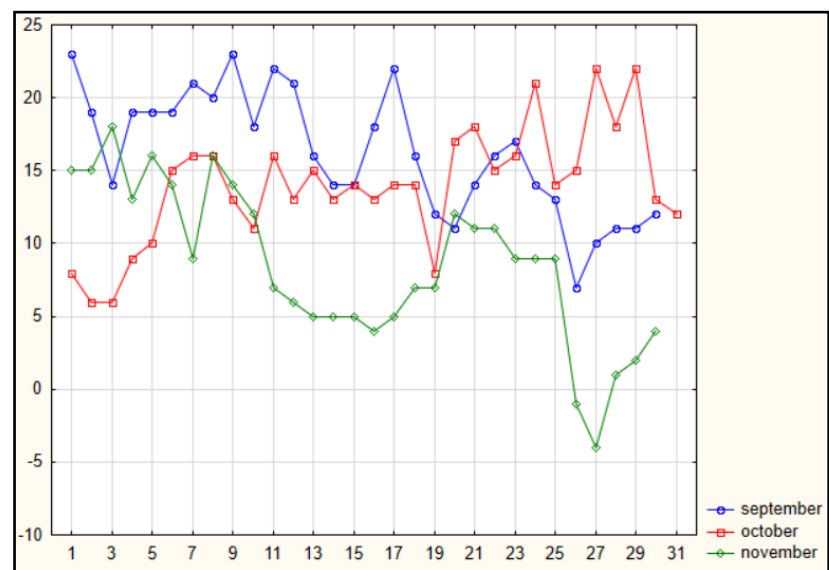


Рис. 8 - Графік температур осінніх місяців

Як бачимо, максимальне нарощення температури відбувається у квітні, а максимальний спад у листопаді. При цьому загалом літні й зимові температури більш нестабільні, ніж осінні та весняні.

Побудуємо графік середньомісячних температур по всьому періоду спостережень (рис. 9).

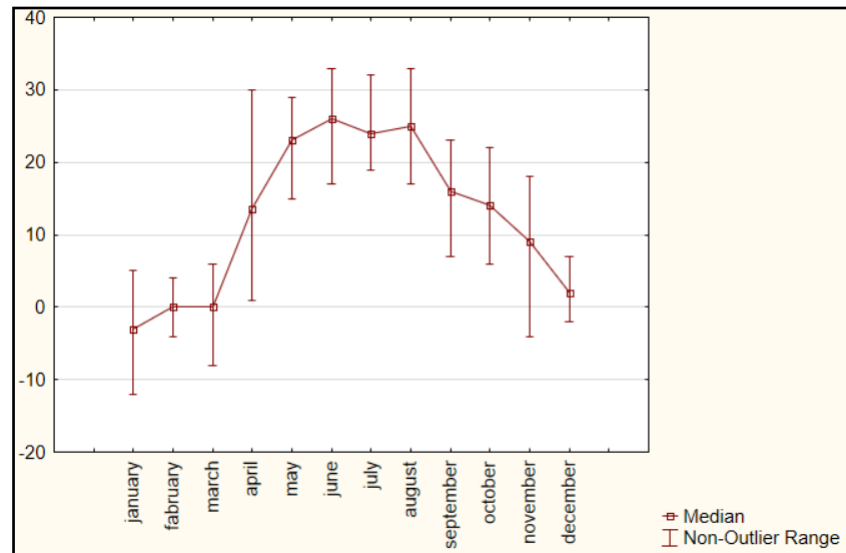


Рис. 9 - Графік середньомісячних температур за 2013 рік

Для виявлення взаємозалежностей у даних скористаємося асоціативними правилами. Для цього перетворимо початковий файл із даними (а точніше - створимо на його основі новий, в якому замість значень температур вкажемо одиниці для тих значень, які більше середнього для відповідного місяця, і нулі для тих, які менше). Щоб це зробити, треба спочатку визначити середні температури для кожного місяця (рис. 10).

Variable	Valid N	Mean	Median	Minimum	Maximum	Std.Dev.
january	31	-2,80645	-3,00000	-12,0000	5,00000	4,245934
fabruary	28	0,25000	0,00000	-4,0000	4,00000	2,187930
march	31	0,06452	0,00000	-8,0000	9,00000	3,492234
april	30	14,06667	13,50000	1,0000	30,00000	8,614076
may	31	22,96774	23,00000	15,0000	29,00000	4,127016
june	30	24,90000	26,00000	17,0000	33,00000	3,594536
july	31	24,87097	24,00000	19,0000	32,00000	3,853500
august	31	25,16129	25,00000	17,0000	33,00000	5,040481
september	30	16,20000	16,00000	7,0000	23,00000	4,286305
october	31	13,96774	14,00000	6,0000	22,00000	4,070085
november	30	8,53333	9,00000	-4,0000	18,00000	5,424932
december	31	2,16129	2,00000	-4,0000	11,00000	2,864690

Рис. 10 - Описові статистики (в т.ч. середні температури місяців - mean)

Отже, новий файл із даними виглядає так:

	1 january	2 fabruary	3 march	4 april	5 may	6 june	7 july	8 august	9 september	10 october	11 november	12 december
1	1	1	1	0	0	0	0	0	1	0	1	1
2	1	1	0	0	0	0	0	1	1	0	1	1
3	1	1	1	0	1	0	1	0	0	0	1	0
4	1	0	0	0	0	0	1	1	1	0	1	0
5	1	1	1	0	0	1	1	1	1	0	1	0
6	1	1	1	0	1	1	1	1	1	1	1	0
7	0	1	1	0	1	1	1	1	1	1	1	0
8	0	0	1	0	1	1	0	1	1	1	1	0
9	0	0	0	0	1	1	1	1	1	0	1	0
10	0	0	0	0	1	1	1	1	1	0	1	0
11	1	0	0	0	1	0	1	0	1	1	0	0
12	0	0	0	0	1	0	0	1	1	0	0	0
13	0	1	1	0	1	0	0	1	0	1	0	0
14	0	1	1	0	0	1	0	0	0	0	0	0
15	1	0	0	0	0	1	0	0	0	1	0	0
16	1	0	0	1	1	1	0	0	1	0	0	0
17	1	0	0	1	1	1	0	1	1	1	0	1
18	0	0	0	1	1	1	1	1	0	1	0	0
19	0	0	1	1	1	0	1	1	0	0	0	0
20	0	0	1	1	1	1	0	1	0	1	1	0
21	1	0	1	0	1	1	0	1	0	1	1	0
22	0	0	0	1	0	1	0	0	0	1	1	1
23	0	1	0	1	0	1	0	0	1	1	1	1
24	0	0	0	1	0	1	0	0	0	1	1	1
25	0	1	0	1	0	1	0	0	0	1	1	0
26	0	1	0	1	0	0	0	0	0	1	0	1
27	0	1	0	1	0	0	1	0	0	1	0	1
28	0	1	0	1	0	1	1	0	0	1	0	1
29	1		1	1	1	0	1	0	0	1	0	1
30	1		1	1	1	0	1	0	0	0	0	1
31	1		1		0		0	0		0		0

Рис. 11 - Новий файл із даними

Тепер знайдемо закономірності між даними за допомогою асоціативних правил (рис. 12).

Min. support = 10,0%, Min. confidence = 50,0%, Min. correlation = 50,0%						
Max. size of body = 10, Max. size of head = 10						
	Body	Head	Support(%)	Confidence(%)	Correlation(%)	
1	june == 0	july == 1	22,58065	58,3333	54,00617	
2	june == 0	august == 0	22,58065	58,3333	50,51815	
3	june == 0	1	32,25806	83,3333	58,92557	
4	june == 0	july == 1, 1	19,35484	50,0000	54,77226	
5	june == 0	august == 0, 1	22,58065	58,3333	50,51815	
6	june == 1	july == 0	35,48387	61,1111	62,88281	
7	june == 1	august == 1	32,25806	55,5556	60,85806	
8	june == 1	1	41,93548	72,2222	62,54628	
9	june == 1	0	58,06452	100,0000	81,64966	
10	june == 1	july == 0, 1	29,03226	50,0000	56,69467	
11	june == 1	july == 0, 0	35,48387	61,1111	66,94387	
12	june == 1	august == 1, 0	32,25806	55,5556	62,99408	
13	june == 1	1, 0	41,93548	72,2222	68,51602	
14	june == 1	july == 0, 1, 0	29,03226	50,0000	61,23724	
15	july == 0	june == 1	35,48387	64,7059	62,88281	
16	july == 0	august == 0	32,25806	58,8235	60,63391	
17	july == 0	1	45,16129	82,3529	69,31033	
18	july == 0	0	48,38710	88,2353	70,01400	
19	july == 0	june == 1, 1	29,03226	52,9412	60,54055	
20	july == 0	june == 1, 0	35,48387	64,7059	62,88281	
21	july == 0	august == 0, 1	32,25806	58,8235	60,63391	
22	july == 0	august == 0, 0	29,03226	52,9412	60,54055	
23	july == 0	1, 0	38,70968	70,5882	65,07914	
24	july == 0	june == 1, 1, 0	29,03226	52,9412	60,54055	
25	july == 0	august == 0, 1, 0	29,03226	52,9412	60,54055	
26	july == 1	june == 0	22,58065	50,0000	54,00617	
27	july == 1	august == 1	25,80645	57,1429	55,20524	
28	july == 1	1	32,25806	71,4286	54,55447	
29	july == 1	0	38,70968	85,7143	61,72134	
30	july == 1	august == 1, 0	25,80645	57,1429	57,14286	
31	august == 0	july == 0	32,25806	62,5000	60,63391	
32	august == 0	1	51,61290	100,0000	81,64966	
33	august == 0	0	41,93548	81,2500	62,54628	
34	august == 0	june == 1, 1	25,80645	50,0000	55,47002	
35	august == 0	july == 0, 1	32,25806	62,5000	66,81531	
36	august == 0	july == 0, 0	29,03226	56,2500	58,09475	
37	august == 0	1, 0	41,93548	81,2500	72,67221	
38	august == 0	june == 1, 1, 0	25,80645	50,0000	55,47002	
39	august == 0	july == 0, 1, 0	29,03226	56,2500	64,95191	
40	august == 1	june == 1	32,25806	66,6667	60,85806	
41	august == 1	july == 1	25,80645	53,3333	55,20524	
42	august == 1	0	45,16129	93,3333	69,56656	
43	august == 1	june == 1, 0	32,25806	66,6667	60,85806	
44	august == 1	july == 1, 0	25,80645	53,3333	59,62848	
56	june == 0, august == 0	july == 1	16,12903	71,4286	50,50763	
57	june == 0, august == 0	1	22,58065	100,0000	54,00617	
58	june == 0, august == 0	july == 1, 1	16,12903	71,4286	59,76143	
59	june == 0, 1	july == 1	19,35484	60,0000	50,70926	
60	june == 0, 1	august == 0	22,58065	70,0000	55,33986	
61	june == 0, 1	july == 1, august == 0	16,12903	50,0000	54,54972	
62	june == 1, july == 0	august == 0	22,58065	63,6364	62,76449	
63	june == 1, july == 0	1	29,03226	81,8182	55,39117	
64	june == 1, july == 0	0	35,48387	100,0000	63,82847	
65	june == 1, july == 0	august == 0, 1	22,58065	63,6364	52,76449	
66	june == 1, july == 0	august == 0, 0	22,58065	63,6364	58,53694	
67	june == 1, july == 0	1, 0	29,03226	81,8182	60,67799	
68	june == 1, july == 0	august == 0, 1, 0	22,58065	63,6364	58,53694	
69	june == 1, july == 1	august == 1	19,35484	85,7143	58,55400	
70	june == 1, july == 1	0	22,58065	100,0000	50,91751	
71	june == 1, july == 1	august == 1, 0	19,35484	85,7143	60,60915	
72	june == 1, august == 0	july == 0	22,58065	87,5000	60,02450	
73	june == 1, august == 0	1	25,80645	100,0000	57,73503	
74	june == 1, august == 0	0	25,80645	100,0000	54,43311	
75	june == 1, august == 0	july == 0, 1	22,58065	87,5000	66,14378	
76	june == 1, august == 0	july == 0, 0	22,58065	87,5000	63,90097	
77	june == 1, august == 0	1, 0	25,80645	100,0000	63,24555	
78	june == 1, august == 0	july == 0, 1, 0	22,58065	87,5000	71,44345	
79	june == 1, august == 1	july == 1	19,35484	60,0000	50,70926	
80	june == 1, august == 1	0	32,25806	100,0000	60,85806	
81	june == 1, august == 1	july == 1, 0	19,35484	60,0000	54,77226	
82	june == 1, 1	july == 0	29,03226	69,2308	60,54055	
83	june == 1, 1	august == 0	25,80645	61,5385	55,47002	
84	june == 1, 1	0	41,93548	100,0000	69,38887	
85	june == 1, 1	july == 0, august == 0	22,58065	53,8462	61,39406	
86	june == 1, 1	july == 0, 0	29,03226	69,2308	64,45034	
87	june == 1, 1	august == 0, 0	25,80645	61,5385	61,53846	
88	june == 1, 1	july == 0, august == 0, 0	22,58065	53,8462	64,71502	
89	june == 1, 0	july == 0	35,48387	61,1111	62,88281	
90	june == 1, 0	august == 1	32,25806	55,5556	60,85806	
91	june == 1, 0	1	41,93548	72,2222	62,54628	
92	june == 1, 0	july == 0, 1	29,03226	50,0000	56,69467	
93	july == 0, august == 0	june == 1	22,58065	70,0000	52,17492	
94	july == 0, august == 0	1	32,25806	100,0000	64,54972	
95	july == 0, august == 0	0	29,03226	90,0000	54,77226	
96	july == 0, august == 0	june == 1, 1	22,58065	70,0000	61,39406	
97	july == 0, august == 0	june == 1, 0	22,58065	70,0000	52,17492	
98	july == 0, august == 0	1, 0	29,03226	90,0000	63,63961	
99	july == 0, august == 0	june == 1, 1, 0	22,58065	70,0000	61,39406	
100	july == 0, 1	june == 1	29,03226	64,2857	56,69467	
101	july == 0, 1	august == 0	32,25806	71,4286	66,81531	
102	july == 0, 1	0	38,70968	85,7143	61,72134	
103	july == 0, 1	june == 1, august == 0	22,58065	50,0000	66,14378	
104	july == 0, 1	june == 1, 0	29,03226	64,2857	56,69467	
105	july == 0, 1	august == 0, 0	29,03226	64,2857	66,71244	
106	july == 0, 1	june == 1, august == 0, 0	22,58065	50,0000	66,14378	
107	july == 0, 0	june == 1	35,48387	73,3333	66,94387	
108	july == 0, 0	august == 0	29,03226	60,0000	58,09475	
109	july == 0, 0	1	38,70968	80,0000	63,24555	
110	july == 0, 0	june == 1, 1	29,03226	60,0000	64,45034	
111	july == 0, 0	august == 0, 1	29,03226	60,0000	58,09475	
112	july == 1, august == 0	june == 0	16,12903	83,3333	58,92557	
113	july == 1, august == 0	1	19,35484	100,0000	50,00000	
114	july == 1, august == 0	june == 0, 1	16,12903	83,3333	64,54972	
115	july == 1, august == 1	june == 1	19,35484	75,0000	50,00000	

Рис. 12 - Результат застосування асоціативних правил

Бачимо, наприклад що якщо температура в липні у певні дні місяця більше середнього, то з великою ймовірністю температура у серпні для відповідних днів також буде більше середнього. Отримані результати - взаємозалежності між змінними - можна згодом використовувати при проведенні кластерного аналізу, побудові регресійних та прогнозових моделей.

КЛАСТЕРНИЙ АНАЛІЗ

Виконаємо кластерний аналіз даних за допомогою методу k-середніх. Нам доведеться заздалегідь задати кількість кластерів (2), які ми хочемо виділити (див. рис. 13).

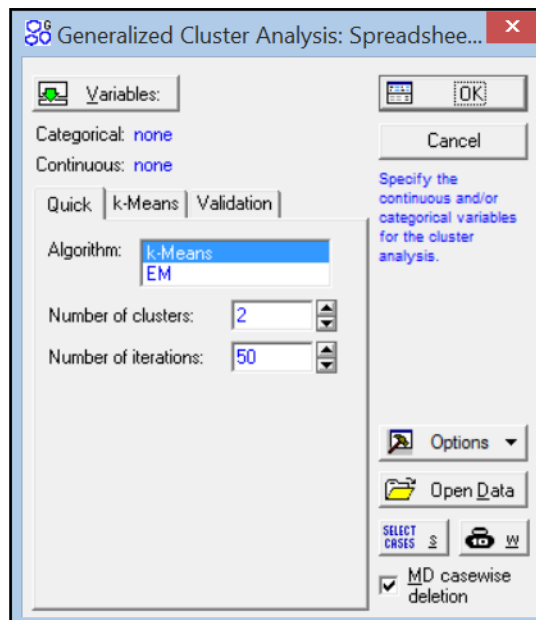


Рис. 13 - Налаштування кластерного аналізу

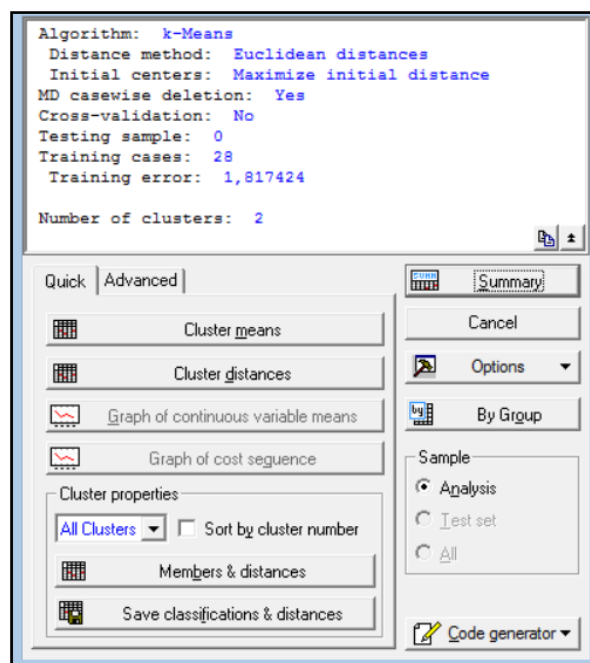


Рис. 14 - Результати аналізу

Number of clusters: 2	
Total number of training cases: 28	
Algorithm	k-Means
Distance method	Euclidean distances
Initial centers	Maximize initial distance
MD casewise deletion	Yes
Cross-validation	No
Testing sample	0
Training cases	28
Training error	1,817424
Number of clusters	2

Рис. 15 - Підсумки по аналізу

		Number of clusters: 2												
		Total number of training cases: 28												
Cluster	january	february	march	april	may	june	july	august	september	october	november	december	Number of cases	Percentage(%)
1	0	1	0	1	0	1	0	0	0	1	0	1	14	50,00000
2	1	0	0	0	1	1	1	1	1	0	1	0	14	50,00000

Рис. 16 - Середні значення змінних для двох знайдених кластерів

Бачимо, що, наприклад, для першого кластеру більш характерні температури серпня менше середнього (15,16), а для другого - більше середнього.

Тепер виконаємо аналогічний аналіз для початкових даних (де не нулі та одинички, а значення температур). Результати цього аналізу представлені на рис. 17-19.

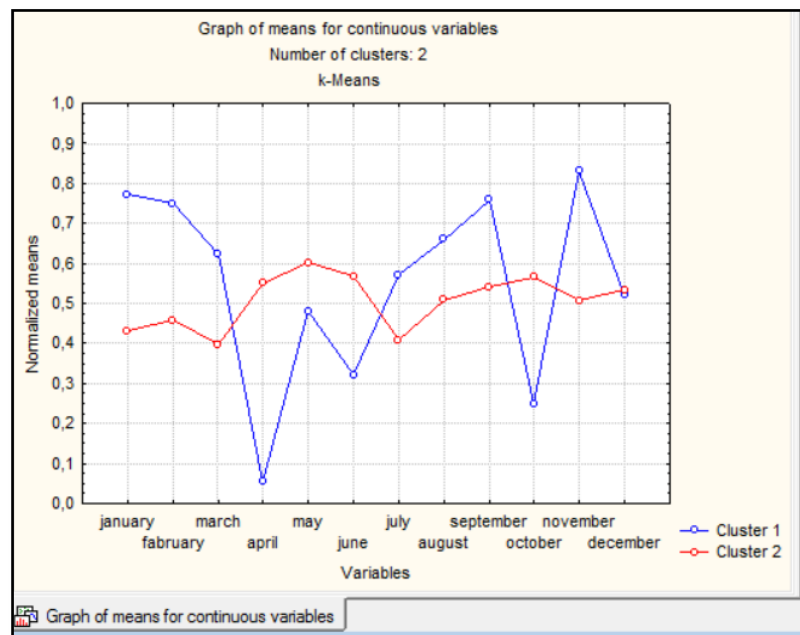


Рис. 17 - Графік нормалізованих середніх значень для двох кластерів

		Number of clusters: 2					
		Total number of training cases: 28					
		Between SS	df	Within SS	df	F	p value
january		469,69	1	307,5238	26	39,711	0,000001
fabruary		30,08	1	100,6667	26	7,770	0,009796
march		78,76	1	243,5238	26	8,409	0,007495
april		6067,98	1	854,6667	26	184,595	0,000000
may		14804,43	1	482,5714	26	797,634	0,000000
june		17707,23	1	272,6667	26	1688,464	0,000000
july		16655,94	1	316,6667	26	1367,540	0,000000
august		18570,96	1	614,2857	26	786,027	0,000000
september		7702,94	1	421,5238	26	475,125	0,000000
october		5441,26	1	292,9524	26	482,921	0,000000
november		2491,07	1	518,0000	26	125,034	0,000000
december		91,18	1	160,0000	26	14,817	0,000692

Рис. 18 - Результати дисперсійного аналізу

Cluster	Number of clusters: 2													
	Total number of training cases: 28													
	january	february	march	april	may	june	july	august	september	october	november	december	Number of cases	Percentage(%)
1	1,14286	2,000000	2,57143	2,57143	21,71429	22,14286	25,85714	27,57143	19,14286	10,00000	14,28571	1,714286	7	25,00000
2	-4,66667	-0,333333	-1,23810	16,95238	23,42857	26,09524	23,90476	25,14286	15,66667	15,04762	7,14286	1,857143	21	75,00000

Рис. 19 - Середні значення температур місяців для двох кластерів

Таким чином, всі місяці мають значний вплив на розділення постережень на кластери; перший кластер містить 7 спостережень (днів), другий - 21; середнє значення температури жовтня для першого кластера становить 14,29, а для другого кластера - 7,14.

ПРОГНОЗУВАННЯ ТЕМПЕРАТУРИ НА 20 ЛЮТОГО 2014 РОКУ

Спочатку спрогнозуємо температуру на всі дні лютого 2013 року. Для цього скористаємося методом Support Vector Machines (методом опорних векторів) (див. рис. 20).

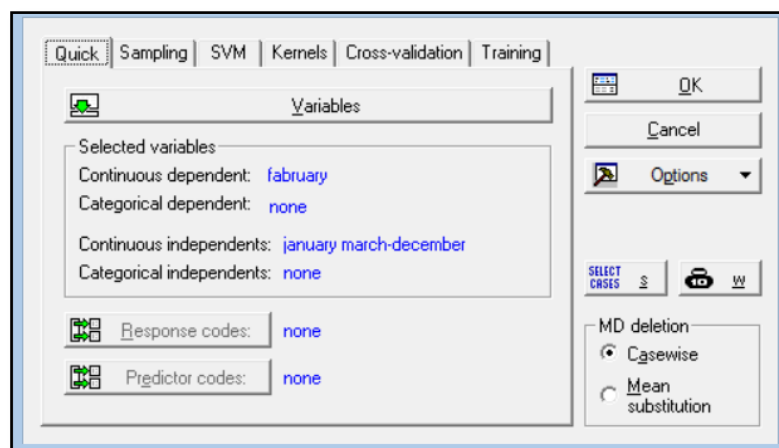


Рис. 20 - Метод Support Vector Machines

Перші 22 спостереження візьмемо для навчання вибірки, інші - для перевірки прогнозів (рис. 21).

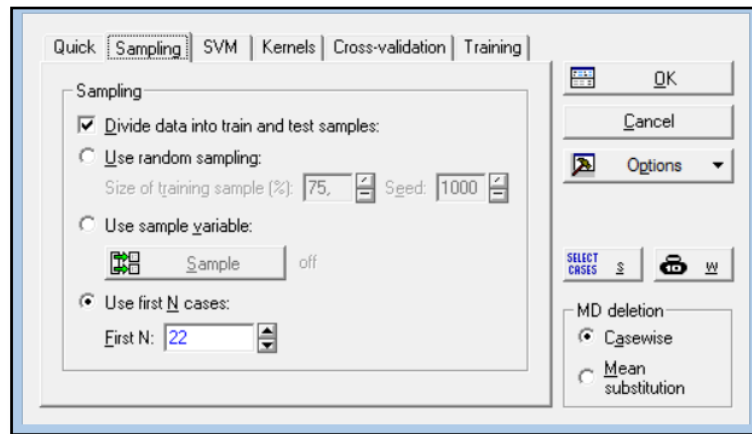


Рис. 21 - Розділення даних на тестові та тренувальні

Результати прогнозування представлені на рис. 22-24.

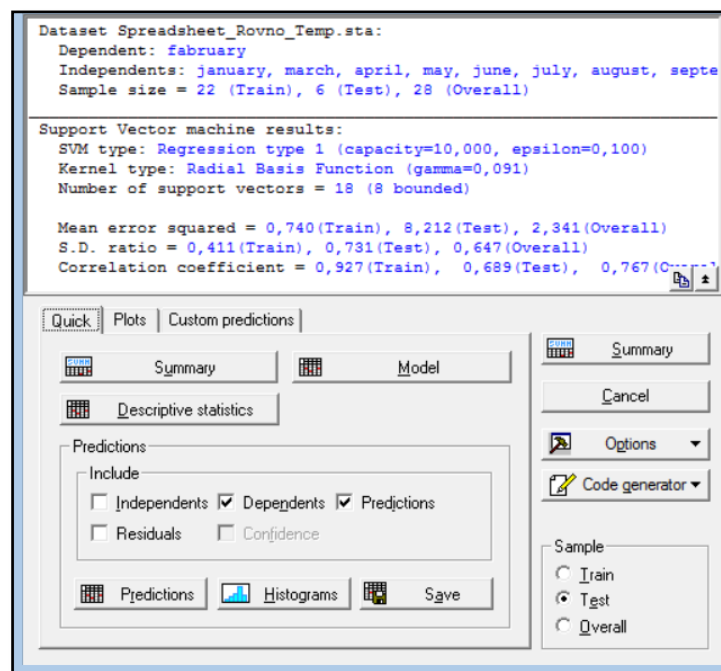


Рис. 22 - Результати прогнозування

		SVM: Regression type 1 (C=10,000, epsilon=0,100), Kernel: Radial Basis Function (gamma=0,091)					
		Number of support vectors= 18 (8 bounded)					
Case Name	fabruary Dependent	fabruary Predicted					
23	1.00000	-2.34163					
24	-1.00000	-1.24857					
25	2.00000	-1.32739					
26	2.00000	-0.23470					
27	4.00000	0.37717					
28	3.00000	0.02500					

Рис. 23 - Результати порівняно з відомими значеннями для тестових спостережень

SVM: Regression type 1 (C=10,000, epsilon=0,100), Kernel: Radial Basis Function (gamma=0,091)			Number of support vectors= 18 (8 bounded)				
Case Name	fabruary Dependent	fabruary Predicted					
1	2.00000	1.60029					
2	1.00000	0.63669					
3	1.00000	0.64487					
4	-1.00000	0.47969					
5	4.00000	1.92724					
6	4.00000	3.78558					
7	3.00000	1.76764					
8	-1.00000	-0.60180					
9	-1.00000	-1.39952					
10	-2.00000	-1.60071					
11	0.00000	-0.40027					
12	-2.00000	-1.58064					
13	1.00000	-0.10760					
14	1.00000	0.60048					
15	-1.00000	-0.59792					
16	-2.00000	-2.77985					
17	-3.00000	-2.60105					
18	-4.00000	-2.00050					
19	-2.00000	-1.68926					
20	-1.00000	-1.40130					
21	-1.00000	-0.47584					
22	0.00000	-0.39821					

Рис. 24 - Результати порівняно з відомими значеннями для тренувальних спостережень