

k-

k-

Modification of K-Means method with unknown number of classes

Isakin Maksim

State University – Higher School of Economics

The classification procedure based on K-Means method with unknown number of classes is designed and investigated in the paper. An adaptive Mahalanobis' metrics and dynamically computed measures of observation anomaly and class homogeneity are used in the method. The classification of regions of the Russian Federation by quality of life integral characteristics is obtained with the proposed method. A brief interpretation of the results is proposed.

1

[, (2006)],

()

[MacQueen (1967)].

，
，
" (. [， (2001)], . 498);
，
，
，
— ，
·
k- ，
· · · · ·
， « » ，
， - ，
， - ，
()
3.
， - ·
，
· · · · ·
： ，
· (. [(2003)]).
-
«
»，

2004-2005 .

$\{S_i\}, i=1,2,\dots,k,$ $k-$ S_i
 $:$ $(p) e_i = (e_{i,1}, e_{i,2}, \dots, e_{i,p})^T$ w_i
 $()$. \ll \gg ,

1. k_0 .

$$k_0^4,$$

$$\begin{cases} e_i = X_i \\ w_i = 1, i = \overline{1, k}. \end{cases}$$

2.

$$d^2(S_i, S_j) = \frac{n_i n_j}{n_i + n_j} (e_i - e_j)^T \Sigma^{-1}(i, j) (e_i - e_j), \quad (1)$$

$$n_k, e_k \quad \Sigma(i, j) - , \quad S_k,$$

$$S_j. \quad \{n_k\} \quad \{\Sigma(i, j)\} \ll \gg^5.$$

4 k_0 ,

5 k_0 .

(

)

$$c_0(S_i, S_j) = \frac{(n_i + n_j - 2)p}{n_i + n_j - p - 1} F_\alpha(p; n_i + n_j - p - 1), \quad (2)$$

$$F_\alpha(v_1, v_2) = 100\alpha - F - v_1 - v_2$$

$$c_0(S_i, S_j), \quad \dots \quad S_i$$

$$S_j \quad \quad \quad S_k$$

$$\begin{cases} e_k = \frac{w_i e_i + w_j e_j}{w_i + w_j} \\ w_k = w_i + w_j \end{cases}$$

,

$$c_0(S_i, S_j). \quad k'_0 < k_0$$

3. X

$$d^2(X, S_j) = \frac{n_j}{n_j + 1} (X - e_j)^T \Sigma^{-1}(j) (X - e_j), \quad (3)$$

S_i .

$$c_1(S_i) = \frac{(n_i - 1)p}{n_i - p} F_\alpha(p; n_i - p). \quad (4)$$

$c_1(S_i)$,

.

,

S_i ,

$$\begin{cases} e_i = \frac{1}{w_i + 1} (w_i e_i + X) \\ w_i = w_i + 1. \end{cases}$$

2,

$$\{n_k\} \quad \{\Sigma(i)\},$$

n_i ,

S_i ,

S_i :

$$\bar{X}(i) = \frac{1}{n_i} \sum_{X_k \in S_i} X_k,$$

S_i :

$$\Sigma(i) = \frac{1}{n_i - 1} \sum_{X_k \in S_j} (X_k - \bar{X}(i))(X_k - \bar{X}(i))^T.$$

$S_i \quad S_j$

$$\Sigma = \frac{1}{n_i + n_j - 2} \left[\sum_{X_k \in S_i} (X_k - \bar{X}(i))(X_k - \bar{X}(i))^T + \sum_{X_k \in S_j} (X_k - \bar{X}(j))(X_k - \bar{X}(j))^T \right].$$

(3)

$\Sigma(i)$ (, ,)

(2) (4)

“ ” (,)

6 . k_0 1% 10% .

($F -$),

, , , .

.

,

,

k_0

.

,

,

.

.

,

- ,

(

)

7 10 (. [

(1981)].

(,

3 7).

,

1, 5 10%

« »

$k-$ (

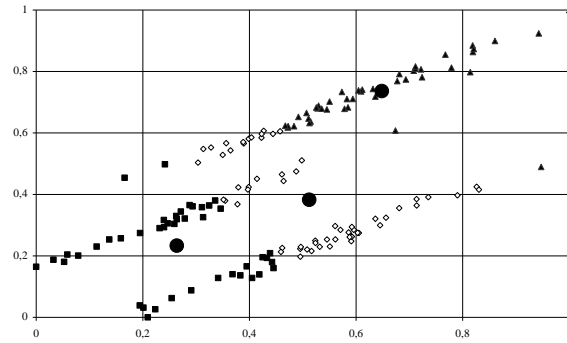
 « »

), ,

 « » (. 1

 $k-$ -

 : —).



. 1.

: ,

 [(2003)]. . 1

 ,

 [(2005)] 2004 .

 , , ,

 , ,

 .

 ,

 .

5%,
2 15.

, 2 -3

1.

	1	2	3	4
1.	315,67	48,65	51,125	97,9
2.				
	312,5	159,21	176,45	247,22
3.				
	20,533	34,363	30,507	21,346
4. 20%	18,633	10,452	10,367	13,646
5.	21,033	18,333	21,435	20,537
6.	9,4333	7,5479	8,5531	10,145
7.	160,73	146,77	131,49	154,6
8.	234,37	87,181	138,84	127,43
9. (1 000)	3,2667	6,2042	2,7227	4,075
1.	61,892	64,35	66,175	63,545
2. 1000	15,831	16,975	12,231	12,061
3.	-7,1979	0,8125	-3,694	-8,4818
4. 100	62,723	28,562	46,149	52,109

000					
5.	100				
000		198,02	162,02	181,89	203,06
6.	100 000	1034,9	549,05	815,01	965,83
7.	100 000	83,379	52,625	61,885	98,376
8.	100 000	61,292	41,075	51,282	69,727
9.	100 000	326,99	239,86	192,07	289,16
10.	1 000	78,183	53,263	66,66	80,745
11.	1 000	1,8188	1,4125	1,7104	1,6364
12.		20,254	26,787	23,21	20,124
13.	()	110,64	550,01	137,96	130,98
14.	23	18,877	17,437	19,612	19,615
1.		8,625	10,247	7,4595	13,073
2.		27,85	28,065	27,576	24,414
3.	1 000 1	0,26071	0,3551	0,35952	0,25405
4.	10 000	12,546	-22,386	-12,105	-10,319
5.	100 000	31,836	26,733	21,181	18,549
6.	100 000	76,671	48,786	44,586	31,397
7.	100 000	9,2536	8,7102	5,1643	4,8054
8.	100 000	406,4	361,19	358,45	269,18
9.	100 000	33,043	50,751	31,514	28,624
10.	,	524,06	210,92	190,15	459,79
11.	,	2110,3	2070,7	3156,6	1768,2

12.	100 000	44,914	60,294	40,2	29,738
13.	,	428,3	62,671	63,752	83,095

«

», «

» «

» .2 .1.

2.

	2	3	4
.	2	1	2
.	3	1	2
.	2	3	4
.	3	3	3
.	3	1	3
.	3	1	3
.	2	3	2
.	3	4	2
.	3	4	3
.	1	2	3
-	4	3	1
.	3	1	2
.	3	1	3
.	3	4	1
-	2	3	4
.	2	4	1
.	3	3	3
.	2	3	2
-	2	3	4
.	3	1	4
.	3	1	2
.	3	1	3
.	2	3	4
.	4	4	4
.	2	4	2
.	3	3	2
.	3	4	1
.	4	3	3

.	4	3	3
.	3	3	1
.	4	3	4
.	3	1	3
.	3	1	3
.	3	3	4
.	3	3	2
.	2	1	1
.	3	1	3
.	3	1	3
.	4	4	3
.	2	3	4
.	3	1	3
.	3	3	4
.	2	1	2
.	3	3	2
.	2	1	2
.	2	3	4
.	3	3	2
.	3	1	3
.	4	4	2
.	3	4	2
.	3	3	2
. ()	4	3	3
.	3	3	4
.	4	3	2
.	2	1	1
.	2	1	2
.	2	3	4
.	3	4	3
.	4	3	1
.	3	1	4
.	4	3	3
.	3	4	1
.	3	1	3
.	2	3	4
.	3	4	3
.	3	1	3
.	3	3	4
.	3	4	3

.	1	2	1
.	2	4	2
.	3	1	1
	3	1	2
.	2	3	1
.	2	1	2
.	3	4	2
. .	1	2	3
.	4	4	3

« »

: . ,

,

.

.

,

(

,

),

,

-

.

,

-

(

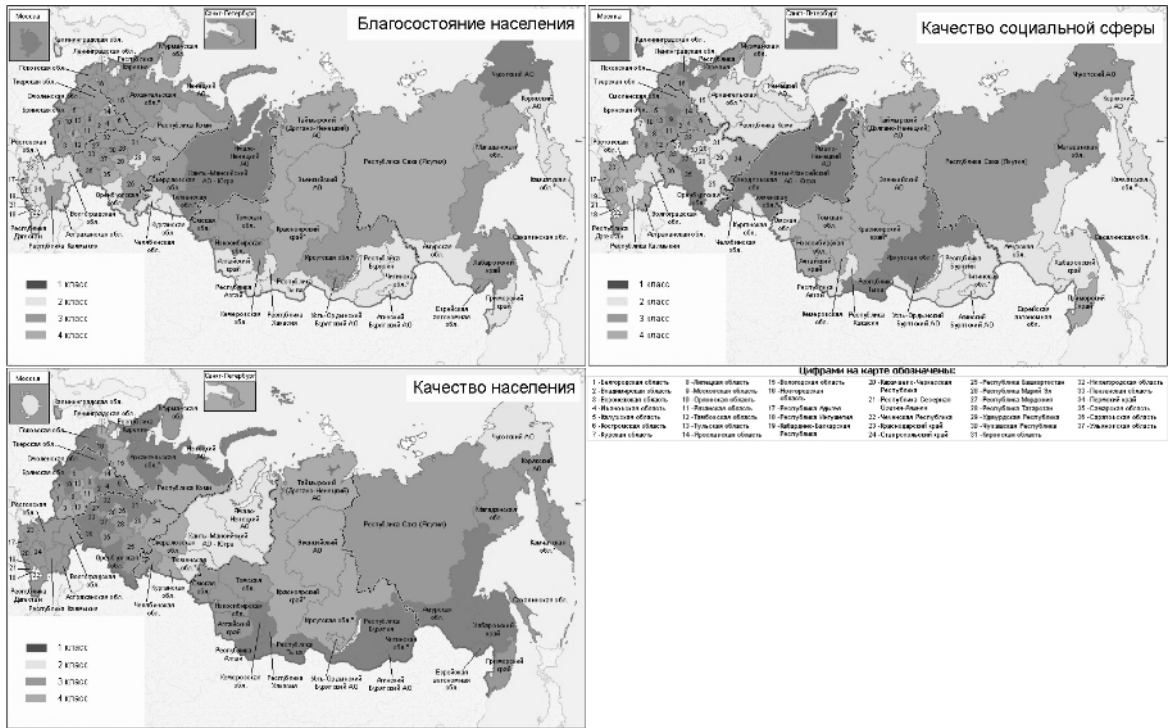
,

).

, ,

,

.



.1.

«

»

«

»

· — , ,
· ()
· « »
· ,
· ,
· ()
·). , ,
· , -
· « »
· « » (),
·
·
· k -
· , · · ·
· « » k - ,
· ,
·

1.
. , 2001.
2.
// ,
39, 2, 2003, . 33-53.
3.
, //
, 2006, 1.
4. //
, " , 1981, 3, . 71 -
77 5, . 153-160.
5. - . - ,
2005.

6. MacQueen, J. Some Methods for Classification and Analysis of Multivariate Observations // Proc. Fifth Berkeley Symp. Math. Stat. and Probab., 1967, vol. 1, pp. 281-297.