

НЕІСРАРХІЧНА КЛАСТЕРИЗАЦІЯ ЗВУКОВИХ ОДИНИЦЬ МОВНОГО СИГНАЛУ

Розглянуто використання кластеризації для класифікації звукових ділянок мовного сигналу на прикладі дифтонгів польської мови. Засобом попередньої мінімаксної нормалізації показана можливість використання методу k-means.

Ключові слова: мовний сигнал, звукові одиниці, дифтонги польської мови, нормалізації, кластеризації, k-means, кореляційна метрика, диз'юнктивне розбиття.

Рассмотрено использование кластеризации для классификации звуковых участков языкового сигнала на примере дифтонгов польского языка. Средством предыдущей минимаксной нормализации показана возможность использования метода k-means.

Ключевые слова: языковый сигнал, звуковые единицы, дифтонги польского языка, нормализации, кластеризации, k-means, корреляционная метрика, дизъюнктивное разбиение.

There is investigated the using of clustering to classify the audio units of speech signal on the example of Polish language diphthongs. The means of the provisional normalization minimax shown to use the method k-means

Key words: speech signal, audio units, Polish language diphthongs, data generalization and clustering, k-means, correlation metric, disjunctive partitioning.

Вступ

Однією із важливих задач попередньої обробки часових рядів є кластеризація даних – поділ їх на відокремлені множини за певними ознаками з подальшим обробленням інформації щодо кожної такої множини відповідно до поставленої задачі. Прикладом кластеризації є класифікація звукових одиниць мовного сигналу в задачах їх аналізу, перетворення та синтезу, оскільки звуки різних класів характеризуються яскраво вираженою індивідуальністю як щодо внутрішньої структури, так і щодо параметричної характеристики їх окремих ділянок. При цьому в залежності від поставленої задачі може змінюватися число класів, на які доцільно поділити звукові одиниці.

У роботі пропонується вирішення задачі класифікації мовних одиниць за тривалістю виділених ділянок окремих звуків. Основою методу є гіпотеза про те, що за вказаними характеристиками мовні одиниці утворюють множини (кластери), які не перетинаються. У зв'язку з цим основою для підтвердження цієї гіпотези пропонується використання методів кластеризації часового ряду, зокрема методу k-середніх (k-means). Завдяки попередній обробці вхідного набору даних вибір методу кластеризації не є принциповим, а зумовлений лише простою практичної реалізації.

Для перевірки якості розробленого методу проведено аналіз даних, що описують внутрішню часову структуру специфічних класів звуків – польських дифтонгів ϱ та φ в різних темпах мовлення в сторону сповільнення. При цьому набір даних включає статистику як твердих дифтонгів, так і м'яких. Ставиться задача встановити, чи сповільнення мовлення допускає (подібно до задачі прискорення темпу) об'єднання усіх дифтонгів в один клас, чи специфіка даної задачі вимагає додаткового поділу за ознакою м'якості.

I. Постановка задачі

Нехай задано скінчений дискретизований на ділянці $[0, T]$ мовний сигнал $x(t)$

$$x(t) = \bigcup_{i=1}^n x_i(t), \quad (1)$$

де $x_i(t)$ – i -та ділянка сигналу $x(t)$, довжина якої визначається тривалістю t_i , n – кількість ділянок. У випадку мовного сигналу ділянка $x_i(t)$ визначає тривалість частини дифтонгу (голосна, приголосна, пауза) в різних темпах.

Відзначимо, що множина $X = \{x_i(t) | i = 1..n\}$ утворює диз'юнктивне покриття [2] розмірності n , тобто має місце умова

$$\forall i, j \in [1, n]: x_i(t) \cap x_j(t) = \emptyset, \quad (2)$$

а тривалості t_i кожної ділянки $x_i(t)$ визначають тривалість T сигналу $x(t)$

$$T = \sum_{i=1}^n t_i . \quad (3)$$

На множині X приймаємо нульову гіпотезу H_0 про те що усі мовні звуки можуть бути класифіковані за характеристиками внутрішньої часової структури ділянок $x_i(t)$ [1]. Тоді метою роботи є перевірка цієї гіпотези.

ІІ. Класифікація як перевірка нульової гіпотези

Для перевірки гіпотези сформує набір багатовимірних векторів в просторі сигналу наступним чином.

За [1] кожну ділянку $x_i(t)$ розб'ємо на три інтервали

$$x_i(t) = \bigcup_{j=1}^3 x_{i,j}(t), \quad (4)$$

які, подібно до випадку сигналу $x(t)$, утворюють диз'юнктивне покриття ділянки $x_i(t)$ і визначають відповідно до значення індексу j стаціонарний, перший та другий переходні інтервали ділянки $x_i(t)$.

Кожен інтервал $x_{i,j}(t)$ є визначенням сигналом $x(t)$ на окремому часовому проміжку тривалості $t_{i,j}$. Сума тривалостей $t_{i,j}$ усіх інтервалів $x_{i,j}(t)$ є рівною тривалості t_i ділянки $x_i(t)$

$$t_i = \sum_{j=1}^3 t_{i,j}, \quad t_{i,j} \leq t_i \leq T . \quad (5)$$

Тоді кожній ділянці $x_i(t)$ у відповідність можна поставити вектор розмірності 3

$$x_i(t) \rightarrow I_i = (t_{i,1}, t_{i,2}, t_{i,3}) . \quad (6)$$

Якщо ділянка $x_i(t)$ є представлена у різних темпах мовлення то це дає можливість розширити розмірність вектора Λ_i

$$\Lambda_i = \{t_{i,j,z} \mid j = 1..3, z = 1..m\} , \quad (7)$$

де m – кількість темпів мовлення, а z – його індекс. Тоді $t_{i,j,z}$ є тривалістю j -го інтервалу $x_{i,j,z}(t)$ i -ої ділянки $x_i^z(t)$ $x_{i,z}(t)$ сигналу $x^z(t)$ в z -му темпі. Розмірність вектора Λ_i рівна $3m$.

В результаті цього мовний сигнал, визначений у різних темпах може бути представлений обмеженим дискретним простором точок \mathbf{R}^{3m} розмірності $3m$

$$x(t) \rightarrow \Xi = \begin{pmatrix} (t_{1,1,1}, t_{1,2,1}, \dots, t_{1,3,m}) \\ \dots \\ (t_{n,1,1}, t_{n,2,1}, \dots, t_{n,3,m}) \end{pmatrix} , \quad (8)$$

який визначається кореляційною метрикою, яка, у свою чергу, у матричній формі має вигляд [2]

$$d(\Lambda_a, \Lambda_b) = 1 - \frac{(\Lambda_a - \bar{\Lambda}_a)(\Lambda_b - \bar{\Lambda}_b)^T}{\sqrt{(\Lambda_a - \bar{\Lambda}_a)(\Lambda_a - \bar{\Lambda}_a)^T} \sqrt{(\Lambda_b - \bar{\Lambda}_b)(\Lambda_b - \bar{\Lambda}_b)^T}} . \quad (9)$$

Тут $\bar{\Lambda}_i = \frac{1}{3m} \sum_{z=1}^m \sum_{j=1}^3 t_{i,j,z}$ – середнє значення елементів вектора Λ_i , а символ T – визначає операцію транспонування.

Набір (8) виступатиме набором вхідних даних і для подальшої перевірки гіпотези H_0 використаємо метод класифікації k -means [3] з метрикою (9).

Оскільки ділянки мовного сигналу, зокрема тривалості $t_{i,j,z}$ інтервалів, характеризуються тим, що вони не розміщаються компактно і не утворюють точкових ущільнень, то безпосереднє використання методу k -means відносно (8) не забезпечить достатньої точності класифікації і, відповідно, підтвердження гіпотези H_0 .

З цією метою усі точки представлення (8) нормалізуємо [4] за схемою мінімаксної нормалізаційної операції. Для цього за (8) сформуємо $3m$ наборів

$$I_k = (t_{1,(k)}, t_{2,(k)}, \dots, t_{n,(k)}) \mid k = (j, z), j = 1..3, z = 1..m . \quad (10)$$

Кожен такий набір I_k розб'ємо на h проміжків $\Delta_{k,q}$ таких, щоб

$$I_k = \bigcup_{q=1}^h \Delta_{k,q}, \quad \forall p, q \in [1; h]: \Delta_{k,q} \cap \Delta_{k,p} = \emptyset . \quad (11)$$

На кожному проміжку $\Delta_{k,q}$ визначимо максимум

$$M_{k,q} = \max_{t_{i,k} \in \Delta_{k,q}} \Delta_{k,q} . \quad (12)$$

Серед усіх максимумів $M_{k,q}$ знаходимо мінімум

$$M_k = \min_{q \in [1; h]} M_{k,q} , \quad (13)$$

обернене значення до якого виступатиме нормалізаційним множником для елементів набору I_k . В результаті за допомогою (13) отримаємо нормалізований набір I'_k

$$I'_k = M_k^{-1} I_k = \left\{ t'_{i,(k)} \mid t'_{i,(k)} = \frac{t_{i,(k)}}{M_k} \right\}_{i=1..n} \quad (14)$$

Якщо набір Ξ представити матрицею розмірності $(n \times m)$, то нормалізаційне представлення буде визначатись добутком ΞN , де N – діагональна нормалізаційна матриця N , яка має такий вид

$$N = \begin{pmatrix} \frac{1}{M_1} & \dots & 0 \\ \dots & \dots & \dots \\ 0 & \dots & \frac{1}{M_{3m}} \end{pmatrix}, \quad (15)$$

У результаті застосування цього представлення отримаємо набір Ξ' , елементами якого є нормалізовані $3m$ -вимірні точки Λ'_i .

Надалі, стосовно набору Λ'_i за метрикою (9) вирішується задача мінімізації дисперсії на усіх точках кожного кластера Ω_s

$$\arg \min \sum_{s=1}^K \sum_{\Lambda_i \in \Omega_s} \|\Lambda_i - \hat{\Delta}_s\|^2, \quad (16)$$

де $\hat{\Delta}_s$ – центр мас (центроїд) кластера Ω_s , K – число кластерів, яке визначається числом класів звуків мови. Сам алгоритм мінімізації визначається методом *k-means*.

За результатами кластеризації проводиться імовірностна оцінка нульової гіпотези H_0 . Фактично рівень значимості тесту визначається сумарним числом помилок в усіх кластерах.

III. Результати практичних експериментів

З метою спрощення обчислювальних затрат розглянемо вхідний набір, який складається лише з твердих та м'яких дифтонгів польської мови. Зокрема у табл. 1, 2 наведено статистики відповідно твердих та м'яких дифтонгів польської мови. У кожному темпі попередньо виділений у мовному сигналі звук був розбитий на три ділянки, тривалості яких наведені у окремому рядку табл. 1 і табл. 2.

Таблиця 1
Тривалості ділянок твердих дифтонгів у різних темпах (клас 1)

1- темп			2- темп			3- темп			4- темп		
t _{1,1,1}	t _{1,2,1}	t _{1,3,1}	t _{1,1,2}	t _{1,2,2}	t _{1,3,2}	t _{1,1,3}	t _{1,2,3}	t _{1,3,3}	t _{1,1,4}	t _{1,2,4}	t _{1,3,4}
190	50	9	184	77	18	255	100	22	351	187	39
181	50	12	243	63	10	288	90	26	373	184	50
154	53	24	205	66	25	214	92	20	381	132	48
157	47	15	190	70	21	203	95	30	329	115	40
165	55	18	170	82	25	183	109	58	183	139	50
178	45	12	179	80	27	196	113	35	206	212	42
159	52	10	155	65	23	249	94	41	213	132	55
164	58	14	172	72	27	176	107	33	208	175	65
175	46	17	194	78	25	243	117	41	229	190	78
155	56	24	181	75	23	220	114	27	312	165	37
160	51	22	188	64	23	231	112	31	337	170	38
163	55	13	201	63	14	236	114	23	309	190	42
172	51	16	188	57	19	227	89	32	278	169	47
177	50	18	215	63	22	246	108	35	346	203	35
168	55	13	207	67	15	226	122	20	297	185	38
159	56	14	185	61	15	213	110	24	263	192	40
145	44	10	191	52	13	221	115	19	260	200	35
166	56	17	189	62	18	218	118	28	237	186	39
162	57	23	200	64	21	230	113	23	269	187	30
156	65	20	189	71	22	215	108	31	245	153	39
159	55	20	199	56	21	231	90	27	255	129	36
161	57	13	193	64	15	228	98	19	274	142	36
178	45	11	234	59	14	262	103	18	312	160	41
155	54	14	177	73	18	193	113	29	249	178	49
170	45	17	197	68	19	201	107	26	264	155	61
145	52	18	167	71	21	188	115	30	244	186	55
159	51	14	170	63	17	190	102	26	253	164	42
163	50	21	196	77	22	225	108	30	300	168	60
172	48	17	231	70	18	249	99	23	330	168	41
173	49	13	216	74	16	225	101	22	280	187	56

Таблиця 2

Тривалості ділянок м'яких дифтонгів у різних темпах (клас 2)

1- темп			2- темп			3- темп			4- темп		
$t_{2,1,1}$	$t_{2,2,1}$	$t_{2,3,1}$	$t_{2,1,2}$	$t_{2,2,2}$	$t_{2,3,2}$	$t_{2,1,3}$	$t_{2,2,3}$	$t_{2,3,3}$	$t_{2,1,4}$	$t_{2,2,4}$	$t_{2,3,4}$
161	49	20	262	60	26	326	84	30	352	108	55
157	60	20	255	65	21	287	79	19	333	107	40
172	53	22	232	55	30	316	84	40	254	186	50
195	50	24	255	70	25	296	83	28	313	177	46
172	51	28	216	60	30	334	80	30	354	156	40
138	60	42	203	70	40	247	107	51	229	200	60
126	45	36	181	74	37	183	110	44	230	171	46
160	52	19	211	70	28	174	113	38	221	161	46
151	47	23	191	68	36	241	98	39	231	194	70
183	49	26	233	53	29	255	88	34	290	165	58
168	46	23	231	57	30	263	91	37	288	181	65
177	52	28	214	67	29	250	81	35	279	156	52
145	48	34	210	63	30	282	87	29	307	142	43
150	51	23	201	57	33	307	77	35	335	152	41
166	46	27	210	54	27	322	84	28	321	135	40
135	52	24	229	57	30	296	101	32	328	182	45
143	59	21	218	64	28	283	97	30	280	160	55
168	47	28	227	65	24	320	91	27	310	148	50
181	62	25	240	54	30	310	75	30	335	127	39
170	56	25	243	73	26	287	98	32	309	175	40
149	52	34	227	63	28	270	88	29	287	155	43
138	44	23	207	66	35	259	96	40	275	143	36
169	47	26	229	57	27	300	83	32	260	160	52
176	51	37	233	50	33	275	78	32	313	141	49
164	46	32	217	60	30	248	94	37	267	162	55
170	49	25	234	63	28	267	82	35	289	147	47
177	55	38	229	60	40	260	75	40	283	156	62
147	45	19	214	57	24	245	98	35	260	159	46
165	47	29	234	65	30	282	88	37	300	155	43
170	49	34	256	70	31	279	101	33	297	188	53

Стосовно цих статистик сформовано два сумарні набори статистик. Один (набір А) із них стосувався статистик лише в одному темпі (дані у табл. 1-2 визначені через позначення «1-темп»), а другий (набір Б) – статистики в усіх темпах. Обидва набори були попередньо нормалізовані за (13)-(15). За нульову приймалась гіпотеза H_0 .

На рис. 1 наведено результати кластеризації за методом $k\text{-means}$ набору А. Чисельні результати цієї кластеризації наведено у табл. 3. В колонці «Клас» табл. 3 через значення 1 та 2 вказано клас вхідної точки, зокрема значення 1 визначає м'який дифтонг, а значення 2 – твердий дифтонг.

За результатами наведеними на рис. 1 та у табл. 3 можна констатувати, що сумарне число помилок становила 15 % від загальної кількості статистик (розмірність вхідного набору рівна 60 точками виду (14)), що дає 85 % рівень значимості прийнятої нульової гіпотези. Такий рівень значимості є достатній для результатів кластеризації, особливо зважаючи на характер вхідних даних.

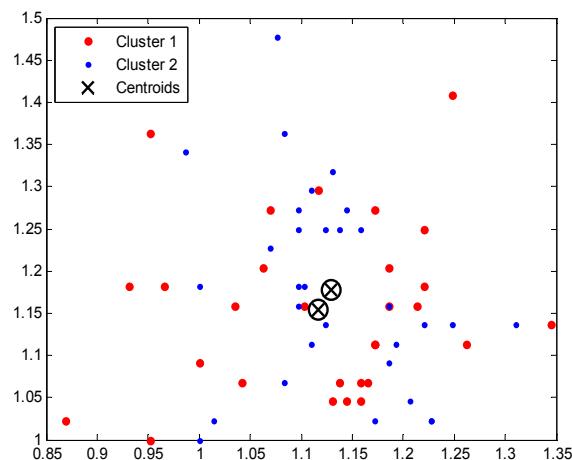


Рис. 1. Візуальне представлення результатів кластеризації попередньо нормалізованого набору А за методом $k\text{-means}$

Таблиця 3

**Чисельні значення результатів кластеризації попередньо нормалізованого набору А
за методом k -means**

Кластер 1				Кластер 2			
Клас	$t_{i,1,1}$	$t_{i,2,1}$	$t_{i,3,1}$	Клас	$t_{i,1,1}$	$t_{i,2,1}$	$t_{i,3,1}$
1	154	53	24	1	190	50	9
1	155	56	24	1	181	50	12
1	160	51	22	1	157	47	15
1	162	57	23	1	165	55	18
2	172	53	22	1	178	45	12
2	195	50	24	1	159	52	10
2	172	51	28	1	164	58	14
2	138	60	42	1	175	46	17
2	126	45	36	1	163	55	13
2	151	47	23	1	172	51	16
2	183	49	26	1	177	50	18
2	168	46	23	1	168	55	13
2	177	52	28	1	159	56	14
2	145	48	34	1	145	44	10
2	150	51	23	1	166	56	17
2	166	46	27	1	156	65	20
2	135	52	24	1	159	55	20
2	168	47	28	1	161	57	13
2	181	62	25	1	178	45	11
2	170	56	25	1	155	54	14
2	140	52	34	1	170	45	17
2	138	44	23	1	145	52	18
2	169	47	26	1	159	51	14
2	176	51	37	1	163	50	21
2	164	46	32	1	172	48	17
2	170	49	25	1	173	49	13
2	177	55	38	2	161	49	20
2	165	47	29	2	157	60	20
2	170	49	34	2	160	52	19
				2	143	59	21
				2	147	45	19

На рис. 2 та у табл. 4. наведено результати кластеризації набору Б. Позначення класу звуку таке ж саме як у випадку результатів, які наведені у табл. 3.

З результатів наведене на рис. 2 та у табл. 4 видно, що число помилок зменшилось до 5 %, що, у свою чергу забезпечує 95 % рівня значимості для нульової гіпотези. Такий високий (зокрема у порівнянні із результатами кластеризації набору А) рівень значимості досягається завдяки збільшенню (у нашому випадку в чотири рази) розмірності вхідних статистик.

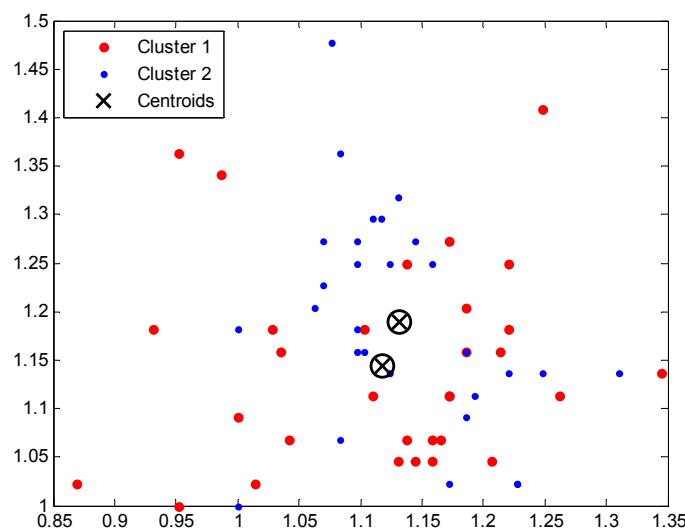


Рис. 2. Візуальне представлення результатів кластеризації попередньо нормалізованого набору Б за методом k -means

Таблиця 4

**Чисельні значення результатів кластеризації попередньо нормалізованого набору Б
за методом *k-means***

Клас	Кластер 1															Клас	Кластер 2														
	Тривалості ділянок																	Тривалості ділянок													
1	165	55	18	170	82	25	183	109	58	183	139	50	1	190	50	9	184	77	18	255	100	22	351	187	39						
1	175	46	17	194	78	25	243	117	41	229	190	78	1	181	50	12	243	63	10	288	90	26	373	184	50						
2	161	49	20	262	60	26	326	84	30	352	108	55	1	154	53	24	205	66	25	214	92	20	381	132	48						
2	172	53	22	232	55	30	316	84	40	254	186	50	1	157	47	15	190	70	21	203	95	30	329	115	40						
2	195	50	24	255	70	25	296	83	28	313	177	46	1	178	45	12	179	80	27	196	113	35	206	212	42						
2	172	51	28	216	60	30	334	80	30	354	156	40	1	159	52	10	155	65	23	249	94	41	213	132	55						
2	138	60	42	203	70	40	247	107	51	229	200	60	1	164	58	14	172	72	27	176	107	33	208	175	65						
2	126	45	36	181	74	37	183	110	44	230	171	46	1	155	56	24	181	75	23	220	114	27	312	165	37						
2	160	52	19	211	70	28	174	113	38	221	161	46	1	160	51	22	188	64	23	231	112	31	337	170	38						
2	151	47	23	191	68	36	241	98	39	231	194	70	1	163	55	13	201	63	14	236	114	23	309	190	42						
2	183	49	26	233	53	29	255	88	34	290	165	58	1	172	51	16	188	57	19	227	89	32	278	169	47						
2	168	46	23	231	57	30	263	91	37	288	181	65	1	177	50	18	215	63	22	246	108	35	346	203	35						
2	177	52	28	214	67	29	250	81	35	279	156	52	1	168	55	13	207	67	15	226	122	20	297	185	38						
2	145	48	34	210	63	30	282	87	29	307	142	43	1	159	56	14	185	61	15	213	110	24	263	192	40						
2	150	51	23	201	57	33	307	77	35	335	152	41	1	145	44	10	191	52	13	221	115	19	260	200	35						
2	166	46	27	210	54	27	322	84	28	321	135	40	1	166	56	17	189	62	18	218	118	28	237	186	39						
2	135	52	24	229	57	30	296	101	32	328	182	45	1	162	57	23	200	64	21	230	113	23	269	187	30						
2	143	59	21	218	64	28	283	97	30	280	160	55	1	156	65	20	189	71	22	215	108	31	245	153	39						
2	168	47	28	227	65	24	320	91	27	310	148	50	1	159	55	20	199	56	21	231	90	27	255	129	36						
2	181	62	25	240	54	30	310	75	30	335	127	39	1	161	57	13	193	64	15	228	98	19	274	142	36						
2	170	56	25	243	73	26	287	98	32	309	175	40	1	178	45	11	234	59	14	262	103	18	312	160	41						
2	149	52	34	227	63	28	270	88	29	287	155	43	1	155	54	14	177	73	18	193	113	29	249	178	49						
2	138	44	23	207	66	35	259	96	40	275	143	36	1	170	45	17	197	68	19	201	107	26	264	155	61						
2	169	47	26	229	57	27	300	83	32	260	160	52	1	145	52	18	167	71	21	188	115	30	244	186	55						
2	176	51	37	233	50	33	275	78	32	313	141	49	1	159	51	14	170	63	17	190	102	26	253	164	42						
2	164	46	32	217	60	30	248	94	37	267	162	55	1	163	50	21	196	77	22	225	108	30	300	168	60						
2	170	49	25	234	63	28	267	82	35	289	147	47	1	172	48	17	231	70	18	249	99	23	330	168	41						
2	177	55	38	229	60	40	260	75	40	283	156	62	1	173	49	13	216	74	16	225	101	22	280	187	56						
2	147	45	19	214	57	24	245	98	35	260	159	46	2	157	60	20	255	65	21	287	79	19	333	107	40						
2	165	47	29	234	65	30	282	88	37	300	155	43																			
2	170	49	34	256	70	31	279	101	33	297	188	53																			

IV. Висновки

Проведені дослідження показали, що при використанні попередньої нормалізації за схемою мінімаксної нормалізаційної матриці існує можливість використання методу *k-means* для даних, які не володіють властивістю компактності і не утворюють точкових згустків (ущільнень) вхідних даних). При цьому зникає потреба у використанні більш складних методів кластеризації та існування розподілів (зокрема нормального як у випадку методу ЕМ) у наборах вхідних даних.

Як свідчать результати експериментів завдяки попередній нормалізації результати роботи алгоритму на основі методу *k-means* у випадку мовних сигналів стали незалежними від вибору базової метрики (різниця результатів в залежності від різних метрик знаходиться в околі одного відсотка).

Розмірність точок кластеризації, яка визначається сумарними статистиками в різних темпах мовлення, безпосередньо впливають на результати кластеризації і з її ростом число помилок кластеризації вже наближається до нуля. Варто зазначати, що статистика навіть одного темпу мовлення в запропонованому методі демонструє достатньо добре результати (рівень значимості –). А сумарні статистики чотирьох темпів мовлення забезпечують практично безпомилкову кластеризацію.

ЛІТЕРАТУРА

1. Рашкевич Ю. М. Перетворення часового масштабу мовних сигналів. – Львів : Академічний експрес, 1997. – 140 с.
2. Новиков Ф. А. Дискретная математика для программистов / Ф. А. Новиков. – 3-е изд. – СПб. : Питер, 2009. – 384 с.
3. Барсегян А. А. Анализ данных и процессов: учеб. Пособие / А. А. Барсегян, М. С. Куприянов, И. И. Холод, М. Д. Тесс, С. И. Елизаров. – 3-е изд. – СПб. : БХВ-Петербург, 2009. – 512 с.
4. Нормализация и распознавание изображений (Публикации международной научно-практической школы «Интеллектуальные системы 2002») [Електронний ресурс] / Е. П. Путятин. – Режим доступу : <http://sumschool.sumdu.edu.ua/is-02/rus/lectures/pyutatin/pyutatin.htm>.