

Харківський національний університет імені В. Н. Каразіна

Факультет математики і інформатики

Кафедра прикладної математики

Кваліфікаційна робота

освітньо-кваліфікаційний рівень: **бакалавр**

на тему

«Математичне моделювання пандемії ковід-19 методами аналізу "великих даних"»

Виконав: студент групи МП41 IV курсу
(перший бакалаврський рівень)
спеціальності 113

«Прикладна математика»

освітньої програми

«Прикладна математика»

Фальченко Іван Русланович

Керівник: доктор фіз.-мат. наук ,

професор кафедри

прикладної математики

Кізілова Наталія Миколаївна

Рецензент: доцент закладу вищої освіти,

кафедри теоретичної та прикладної

інформатики,

кандидат технічних наук

Меняйлов Євген Сергійович

Харків – 2024 рік

Зміст:

1. Анотації.....	3
2. Вступ.....	5
3. Аналіз літератури	7
4. Матеріали та методи.....	11
5. Результати.....	19
6. Висновки.....	30
7. Список використаних джерел.....	32
8. Додатки.....	37

1. Анотації

Фальченко І. Р. Математичне моделювання пандемії ковід-19 методами аналізу "великих даних"

В основі цієї роботи лежить кластеризація європейських країн на основі даних про захворюваність на COVID-19 з 2020 по 2022 роки, які містять кількість нових випадків на мільйон населення та загальну кількість випадків на мільйон населення. Основними методами кластеризації, які використовуються в дослідженні, є модифікації метода k-середніх на основі евклідової метрики та відстані Махаланобіса, які були розглянуті як швидша альтернатива до методу найближчого сусіда. Через особливості методу k-середніх, цей алгоритм було застосовано до 2 наборів початкових центроїдів: 4 початкових центроїдів, які були побудовані на основі надійних результатів кластеризації того ж самого набору даних за допомогою методу найближчого сусіда з дослідження [18] та 7 початкових центроїдів на основі географічних підгруп європейських країн. Метод найближчого сусіда був додатково застосований до кластерів, виявлених за допомогою кластеризації k-середніх, щоб дослідити внутрішню структуру цих кластерів. Результати підтвердили важливість географічної близькості країн, виділяючи при цьому незвичайні сполуки географічно віддалених країн, які неодноразово з'являлися в різних результатах кластеризації.

Falchenko I. R. Mathematical modeling of the covid-19 pandemic using big data analysis methods

This paper is dedicated to the clusterization of European countries based on their COVID-19 morbidity data from 2020 to 2022, which include the number of new cases per million population as well as the total number of cases per million population. The main methods of clusterization used in the research are the variations of k-means clustering based on Euclidean metrics and Mahalanobis distance, which were tested as a faster alternative to the nearest neighbor method. Due to the special features of k-means clustering this algorithm was applied to 2 sets of initial centroids: 4 initial centroids that were constructed based on the reliable results of clusterization of the same dataset using the nearest neighbor method from the research [18] and 7 initial centroids based on the geographical subgroups of European countries. The nearest neighbor method was additionally applied to the clusters revealed by k-means clustering in order to investigate the inner structure of these clusters. The results confirmed importance of geographic vicinity proximity while showing the unexpected combinations of geographically distant countries that repeatedly appeared in different clusterizations.

2. Вступ

В останні часи проблема пандемії COVID-19 нібито втрачає актуальність в суспільній думці, але попри зростання ролі інших глобальних проблем, це питання все ще залишається нагальним. Протягом 2023 року швидкість розповсюдження коронавірусу в деяких країнах Європи досягала більш ніж 500 нових випадків захворювання на день, і цей показник є вочевидь далеким від прийняттого мінімуму. Ця статистика демонструє, що фінальне рішення все ще не знайдене і наукова спільнота має продовжувати роботу в цьому напрямку. [41]

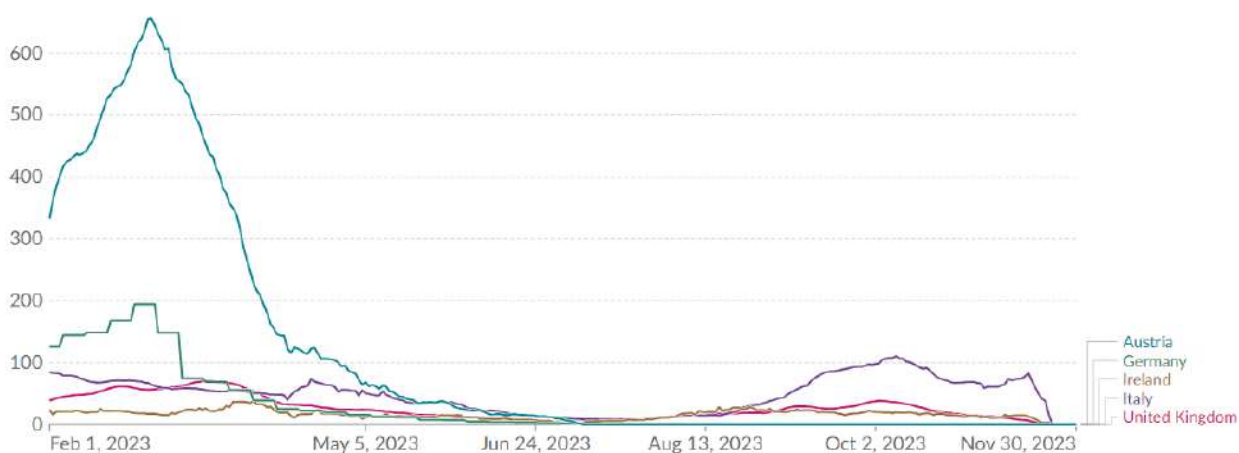


Рис. 1 Графіки нових випадків захворювання на COVID-19 на мільйон осіб для Австрії, Німеччини, Ірландії, Італії та Великобританії з 8 січня 2023 р. по сьогоднішній день. Дані взяті з джерела [32]

COVID-19 – це висококонтагіозна вірусна хвороба, спричинена коронавірусом гострого респіраторного синдрому SARS-CoV-2. Перші випадки цього захворювання були зареєстровані наприкінці 2019 року в місті Ухань провінції Хубей (Китай). Згодом, наприкінці грудня 2019 року, SARS-CoV-2 швидко поширився по всьому світу і, згідно з останніми оновленими статистичними даними, виміряними в 231 країні та країні, він спричинив понад 7

010 681 смертей і 675 619 811 людей одужало. У 2020-2021 роках були виявлені наступні нові варіанти вірусу [4]: Альфа у Великобританії (УК) наприкінці грудня 2020 року; Бета в Південній Африці в грудні 2020 року; Гамма в Бразилії на початку січня 2021 року; Дельта в Індії в грудні 2020 року; Лямбда в Перу в червні 2021 року; і Омикрон в Південній Африці в листопаді 2021 року. Крім того, ще кілька варіантів були класифіковані як Епсилон, Дзета, Ета, Тета, Йота, Каппа і Мю [4]. Ці спостереження демонструють, що географічний фактор має істотне значення у появі нових осередків хвороби та її поширенні навколо них. [41]

В рамках дослідження пандемії COVID-19 дуже важливу роль відіграє аналіз статистичних даних, таких як кількість нових випадків захворювання, нових смертей, нових тестів та щеплення, наявність лікарів та лікарень, коефіцієнт репродукції хвороби (він обчислюється як кількість вторинних випадків зараження однією інфікованою особою) та деяких інших параметрів, що впливають на прийняття рішень щодо заходів протидії пандемії в різних країнах і регіонах [1,3]. В залежності від типу аналізу можна отримати інформацію про ймовірні причини, впливові фактори, схожість чи відмінність ситуації в окремих групах населення та оцінити ефективність уже застосованих заходів [5]. Різні статистичні підходи на основі багатовимірних лінійних та нелінійних регресійних моделей, моделей нелінійної цілочисельної регресії, моделей глибоких нейронних мереж та інших моделей машинного навчання були перевірені на відкритих наборах даних (ВООЗ, статистика коронавірусу Worldometer, Google Health та деякі інші). Математичне моделювання динаміки пандемії [6-8], аналіз стабільності та ефективне керування [9,10] разом з розробкою нових вакцин [11] є важливими для подальшого розвитку математичної епідеміології [12] як галузі міждисциплінарних досліджень. [40, 41]

У цьому дослідженні кластерний аналіз, що являє собою швидкий та надійний математичний інструмент для вирішення задач розпізнавання та

класифікації, був застосований до наборів даних про випадки захворюваності на COVID-19 у 2020, 2021 та 2022 роках. Кластерний аналіз, як частина цього розвиненого наукового апарату, досліджує саме ступінь подібності даних і дозволяє побудувати класифікацію в рамках різноманітних соціальних груп (наприклад, населення різних країн або статей). Результати кластерного аналізу надають нам системне уявлення про ситуацію та дозволяють робити обґрунтовані припущення щодо впливу інших показників на динаміку захворюваності при одночасному порівнянні схожості за статистичними даними та схожості за значеннями цих показників. [40, 41]

В рамках кластерного аналізу даних про коронавірус природнім чином встає питання вибору між різними методами. Кожного разу виникає потреба знайти метод, що є оптимальним за швидкістю, точністю та наочністю результатів для конкретної задачі. Одним з найпростіших засобів порівняння різних методів є порівняння їх роботи на однакових даних. В цьому дослідженні були порівняні результати роботи методу найближчого сусіда та методу k-середніх завдяки застосуванню цих методів до даних про захворюваність на COVID-19 в країнах Європи. Задля дослідження особливостей роботи методу k-середніх він також був застосований до кумулятивних даних з джерела [32] (тобто даних про загальну кількість хворих на COVID-19 в країнах Європи) та була розглянута модифікація цього метода з відстанню Махаланобіса. [40, 41]

3. Аналіз літератури

Загальна описова статистика [5,13], кореляція динаміки захворюваності на covid-19 із соціально-економічними факторами [6], вплив віку та статі [7] та часова динаміка передачі COVID-19 [14,15] з 2020 по 2023 рік були проаналізовані за допомогою різних статистичних методів. Також були

опубліковані локальні дослідження для населення Південно-Східної Азії та Індії [3], Бразилії [16], України [17], Європи [18,19], Польщі [10], Перу [20], Індонезії [21], Китаю [22], США [23] та Канади [24]. Статистичний аналіз і математичні моделі виявились дуже ефективними в прогнозуванні динаміки пандемії у різних країнах [22,23], соціально-економічних факторів [6,25] і урядових заходів проти COVID-19 [26]. Найбільш перспективні результати були отримані на основі світових даних із більш ніж 180 країн [25,27,28]. Підхід до багатовимірного аналізу даних [29], імовірнісні та стохастичні методи [16,30] і глибоке машинне навчання [5] визнано найбільш перспективними методами аналізу великих даних пандемії та ефективною ідентифікації кластерів COVID-19 у світі [28]. Серед усього різноманіття статистичних методів кластерний аналіз є важливим засобом рішення задач класифікації, особливо в випадках, коли щільність ймовірності для вихідних даних (що зазвичай подані в вигляді часових рядів $\{X(t_j)\}_{j=1}^N$, де N – це кількість спостережень, тобто днів) значно відрізняється від нормального розподілу і тому кореляційний аналіз може дати помилкові результати в задачах розпізнавання та класифікації [31]. [41]

Застосування методу k-середніх в цій роботі буде спиратися на результати дослідження «Аналіз «великих даних» і математичне моделювання епідемії covid-19 в країнах Європи» (джерело [18]), а саме на дерево кластеризації, що було отримано за допомогою методу найближчого сусіда. Як буде показано далі, це рішення є цілком природнім, оскільки метод k-середніх вимагає задання початкових центроїдів, які було б доречно визначити за допомогою вже існуючих результатів кластеризації обраних даних.

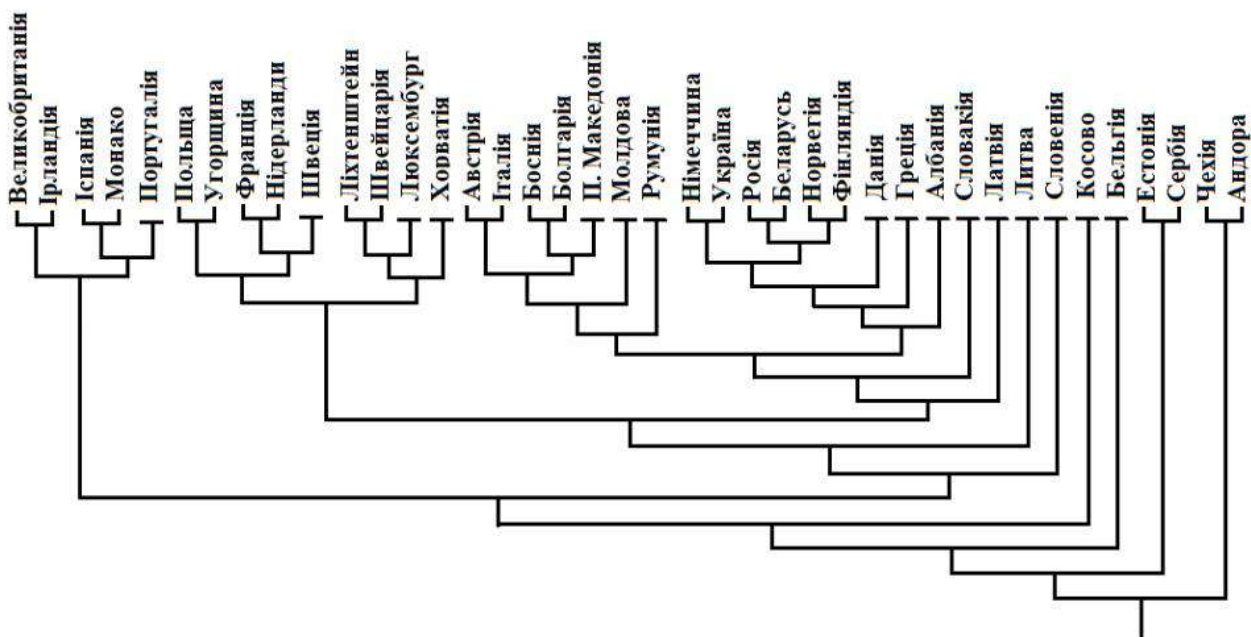


Рис. 2 Дендрограма, що ілюструє результати кластеризації даних про захворюваність в країнах Європи за допомогою методу найближчого сусіда (взято з джерела [18])

У дослідженні [18] кластери для динаміки COVID-19 з березня 2020 року до жовтня 2021 року у 40 європейських країнах були визначені методом найближчого сусіда та візуалізовані у вигляді дендрограм, що показують результати ієрархічної кластеризації (Рис.2). Було помічено, що багато країн з найближчих кластерів мають спільні географічні кордони, наприклад Великобританія та Ірландія, Іспанія та Португалія, Австрія та Італія, Росія та Білорусь, Норвегія та Фінляндія (рис.2). Однак деякі з найближчих сусідів в динаміці COVID-19 знаходяться далеко один від одного на мапі. Цю динаміку можна пояснити спільними заходами уряду проти COVID-19 або іншими соціально-економічними факторами у різних країнах. Також важливо відзначити, що застосування методу найближчого сусіда не є доцільним для великої кількості країн (більш ніж 200), кожній з яких відповідають досить довгі часові ряди

(зазвичай з кількохсот днів спостереження), оскільки алгоритм має досить велику обчислювальну складність (щонайменше $O(n^2)$ в залежності від реалізації, де n – кількість країн). [41]

Неоднозначність результатів разом з занадто малою швидкістю кластеризації можна виправити, якщо застосувати до тих самих даних інший алгоритм кластеризації, наприклад метод k-середніх. Метод k-середніх вже був успішно застосований для класифікації випадків COVID-19 на основі неповних гетерогенних наборів даних [32,33] разом з поточними параметрами атмосфери [34], а також як метод виявлення COVID-19 за послідовностями генома людини [35]. Окрім цього, метод k-середніх продемонстрував свою ефективність в рамках класифікації за допомогою моделей машинного навчання [33]. [41]

На дендрограмі з дослідження [18] можна виділити 4 кластери, за якими буде розрахований один з наборів початкових центрів для методу k-середніх.

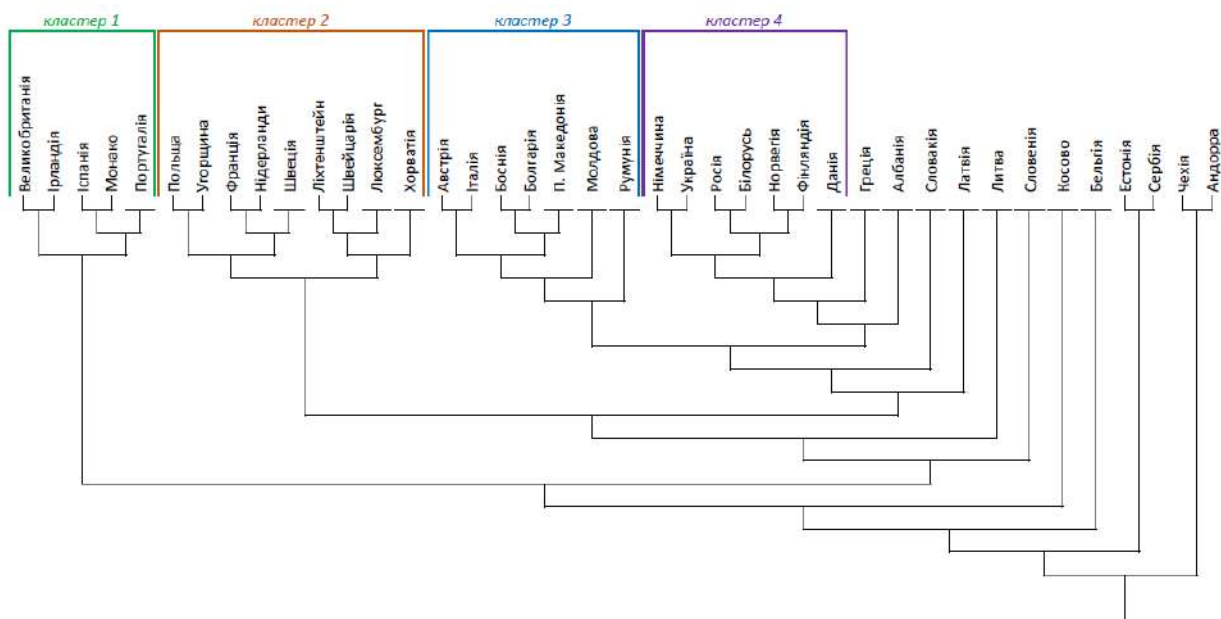


Рис. 3 Дендрограма з рис. 2 з виділеними кластерами

Треба також відзначити, що існуюча наукова практика показує ефективність і точність відстані Махаланобіса в багатьох задачах класифікації, кластеризації та виявлення викидів (включаючи кластеризацію даних COVID-19), оскільки вона враховує кореляції набору даних і тому її використання призводить до більш коректних результатів, особливо коли характеристики в наборі даних корельовані або коли набір даних в цілому є зміщеним та дуже нерівномірним. [33]. Ці переваги відстані Махаланобіса будуть використані в подальшому задля того, щоб скорегувати попередні результати кластеризації, які були отримані при застосуванні евклідової метрики. [41]

4. Матеріали та методи

На першому етапі були використані статистичні дані про кількість нових випадків захворюваності на COVID-19 на мільйон населення в країнах Європи з 14 березня 2020 р. до 15 жовтня 2021 р. (ці дані були взяті з джерела [32] та використані в дослідженні [18] для побудови дендрограми на рис. 2), що подані у вигляді часових рядів (часовий ряд – це послідовність спостережень, які впорядковані за часом у вигляді масиву точок або неперервної функції). Для роботи з ними було використане програмне забезпечення Excel.

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
1	date	Europe	Albania	Andorra	Austria	Belarus	Belgium	Bosnia	Bulgaria	Croatia	Czech	Denmark	Estonia	Finland	France	Germany	Greece	Hungary	Ireland	Italy
2	2020-03-14	12,063040	2,6260158	68,822241	20,874808	0,4671075	14,44098	1,762678	1,4445929	2,1628571	10,842264	17,95212	6,573823	12,05973	16,238063	3,948101	0,6997436	0,2778913	0,1940037	
3	2020-03-15	14,502339	2,1707159	90,040798	25,023577	0,5122493	17,153558	2,192037	1,674432	3,4427448	6,9207359	18,854157	15,16883	7,2875688	15,787425	19,517161	3,8207808	0,838024	11,827659	57,193517
4	2020-03-16	16,265285	2,3366447	114,05405	29,700643	0,9509892	20,322023	2,811099	1,9747762	4,337551	7,7578214	18,769003	20,01289	8,013029	18,125848	22,823187	4,082327	1,003871	14,289250	62,417885
5	2020-03-17	19,779234	2,529588	140,24886	34,886988	0,9777153	23,974703	3,1988885	2,151029	5,3881735	9,9084169	18,782993	20,65314	8,7825888	30,552537	26,123204	4,3316092	1,199117	17,26675	67,489764
6	2020-03-18	21,388872	2,7696376	167,80013	40,86283	0,9888363	28,194304	3,8462718	2,371945	6,8233873	10,284703	18,782775	21,168983	9,843299	22,384979	29,788286	4,9201451	1,6187163	20,478495	72,273285
7	2020-03-19	25,511121	3,1660915	196,44434	46,238838	0,8032551	33,653695	4,208514	2,9855294	7,7091837	11,6444	17,091124	21,244805	10,381705	26,257275	33,500599	4,183844	1,6632929	23,185604	76,632379
8	2020-03-20	28,01252	3,4094885	226,47821	52,09281	0,6021824	38,68278	4,7802211	2,7757143	9,0968367	13,041101	17,888321	22,578822	11,192884	29,359405	37,309384	5,1189205	2,27194	88,43027	
9	2020-03-21	28,558209	3,7847226	254,07193	57,77978	0,6095384	44,588714	5,4025439	2,9430361	10,375	14,423177	18,589363	23,769423	12,82558	32,666927	41,141269	5,428421	2,222241	30,777327	83,548205
10	2020-03-22	31,05333	4,1764519	281,25185	62,980148	0,6307485	51,864886	6,0682598	3,0873344	11,61032	15,018377	19,82331	25,78347	12,822647	36,128107	44,978239	5,780218	2,259479	39,268995	89,910284
11	2020-03-23	33,354452	4,5578043	305,21769	67,37281	0,6888037	59,359684	6,7888088	3,2081125	12,750489	17,344888	21,383084	27,287877	13,547544	39,358891	48,713882	6,1038867	2,6354248	37,188391	87,589785
12	2020-03-24	35,865826	4,9188388	328,68895	70,73027	0,8985439	67,157421	7,5588758	3,3695752	13,828873	18,71816	23,232222	28,419037	14,182889	43,178622	52,254222	6,4521153	3,1402381	41,018422	88,343268
13	2020-03-25	37,954984	5,222334	343,04127	72,919111	1,002768	74,988007	8,3828842	3,598826	14,70694	19,966103	25,326554	31,828526	14,72525	46,584102	55,583345	7,3331117	3,4331117	44,218325	88,488274
14	2020-03-26	39,755825	5,483287	354,31588	73,845321	1,3043185	82,482216	9,2498884	3,6788506	15,622245	21,138772	27,124355	33,757508	15,170666	52,899192	58,377842	7,1183411	3,7038501	47,223559	87,954387
15	2020-03-27	41,201595	5,6968771	368,58853	73,39336	1,7430085	89,419803	10,174600	3,439346	16,948796	22,8992	30,315419	35,647181	15,58851	61,481443	68,77829	7,3988858	3,9491928	50,246882	88,897388
16	2020-03-28	42,245873	5,8768142	361,04077	71,675238	2,3468438	95,563477	11,06883	3,2802305	16,285408	22,859825	33,096017	37,26573	16,038886	70,370769	62,599199	7,8152662	4,968843	53,12382	85,414084
17	2020-03-29	42,855104	6,0358296	357,30947	68,888843	3,1530688	100,72954	11,807858	3,6437888	16,45102	22,40881	30,983421	36,570891	16,588865	68,584307	63,798785	7,789234	5,001291	50,00129	83,69371
18	2020-03-30	43,026983	6,1784573	350,37435	65,303814	4,1863293	104,89338	12,845895	3,6788308	18,495812	23,745778	30,852797	36,584113	17,22817	72,28817	64,388223	7,819686	4,5517586	50,029485	81,823516
19	2020-03-31	42,945449	6,2888081	341,24285	61,14258	5,5017538	108,07803	13,284577	3,8827538	19,342867	23,888803	31,592355	40,088576	18,048441	55,54584	64,377522	7,792796	4,908196	67,28355	78,524211
20	2020-04-01	42,588073	6,3691145	331,18506	56,671581	7,1063964	118,54328	13,752134	3,887384	19,128285	23,817825	31,075277	40,259344	19,028605	61,148851	63,83317	7,8889883	4,9818189	68,20015	77,408413
21	2020-04-02	42,116573	6,4889282	320,83654	52,117056	9,8054273	112,37405	14,118152	3,8991409	15,02888	20,586655	28,182447	40,011159	20,127663	224,27465	62,791756	7,3004789	5,268245	70,88021	75,349044
22	2020-04-03	41,66901	6,599016	310,43417	47,781758	11,268881	113,77285	14,322336	3,61111	15,828939	23,181051	27,890286	36,368481	21,257088	247,18883	61,311845	7,242875	5,022526	70,59179	73,387028
23	2020-04-04	41,320895	6,7188188	300,33332	43,729864	13,821511	114,87818	14,484293	3,5480247	15,282065	22,612935	28,454708	36,328132	22,338221	268,54091	59,454889	6,393569	6,056374	63,440983	71,574825
24	2020-04-05	41,378888	6,7819302	290,17325	40,07392	16,039895	115,78849	14,542741	3,4823553	14,802347	21,87701	28,698794	36,851844	23,224239	265,19282	57,31627	6,8009175	6,559381	59,08482	69,56044
25	2020-04-06	41,125881	6,8982811	278,89511	36,793816	19,852111	118,54837	14,594863	3,4172191	14,702445	20,91789	28,184895	36,027861	23,832788	261,7837	54,938245	6,2841786	7,0881254	60,88395	68,442864
26	2020-04-07	42,388740	5,7986851	258,58885	33,873758	22,798388	117,20138	14,906652	3,368682	14,474388	19,978793	26,953868	33,138004	24,175486	249,59871	52,372017	5,877885	7,6243146	116,89581	67,019011
27	2020-04-08	43,181711	5,6029482	258,53789	31,244123	26,044472	117,58824	14,402888	3,242247	14,248102	18,054504	24,228888	31,151432	24,177801	228,11917	48,673986	6,4848183	6,117821	121,26311	65,71362
28	2020-04-09	44,323804	5,4408884	243,5728	28,822099	29,191849	117,84886	14,376681	3,2888131	13,985571	17,058864	23,118664	28,225448	23,887988	204,51025	46,912006	5,11883	6,533886	131,62384	64,442588
29	2020-04-10	45,429135	5,2438631	229,9343	26,518831	32,901881	117,87388	14,288571	3,241983	13,854449	16,486821	20,403207	27,470754	23,38489	166,00822	44,168254	4,735854	6,843213	140,59182	63,188085
30	2020-04-11	46,302917	5,069332	215,84372	24,39772	35,227708	117,62895	14,008831	3,2414786	13,222143	15,362187	18,589786	25,351466	22,799467	171,1097	41,529384	4,3057697	6,6291587	147,8392	61,896398
31	2020-04-12	46,302917	5,069332	207,64546	22,188281	38,042803	117,88213	13,324481	3,2112591	12,87889	14,288325	18,453451	24,882886	22,222786	158,29148	38,027454	4,803818	6,0928889	153,11773	68,566279
32	2020-04-13	47,07403	5,0947065	197,14885	23,00047	40,72883	118,18888	13,536227	3,0148884	12,066429	13,325903	14,594852	20,02898	21,729482	148,0192	36,701341	3,6788888	6,0718227	156,88881	68,169423

Рис. 4 Частина вихідних статистичних даних з джерела [32]

До вихідних даних можна також віднести дендрограму на рис. 2, оскільки на її основі ми будемо задавати один з наборів початкових центроїдів для методу k-середніх.

Тепер перейдемо до основної складової цієї роботи, а саме кластерного аналізу. Головне призначення кластерного аналізу - розбивка безлічі досліджуваних об'єктів і ознак на однорідні у відповідному розумінні групи або кластер, що є задачею класифікації даних і виявлення в них структури в ній.

Задача кластерного аналізу полягає в тому, щоби різномірні дані, які містяться в множині Q , розбити на ціле число m підмножин (кластерів) Q_1, \dots, Q_m так, щоб кожен об'єкт G_j належав до однієї і тільки однієї підмножині розбиття, і щоб об'єкти, які належать одному і тому ж кластеру, були подібними, в той час як об'єкти, які належать різним кластерам, були різномірними. [40]

Загальна схема кластерного аналізу на першому етапі має наступний вигляд:

I. Проводимо кластеризацію методом k-середніх. В якості результату отримуємо набір кластерів та їх центри. Для подальшого аналізу обчислюємо радіуси кластерів та відстані від кожного елемента до центра того кластера, якому він належить.

II. Оскільки результати застосування методу k-середніх та методу найближчого сусіда мають різний вигляд, треба застосувати деякі перетворення, так щоб подати кластери, що були отримані методом k-середніх, у вигляді дендрограми. Для цього визначимо нову відстань між елементами, спираючись на відстані елементів до центру відповідного кластера, та застосуємо метод найближчого сусіда з цією метрикою всередині кожного кластера. Таким чином

отримуємо окремі дерева для кожного кластера, які можна буде порівняти з кластерами на дендрограмі на рис. 3.

Розглянемо основні характеристики та особливості методів кластеризації, що будуть застосовані у подальшому:

1) Метод k-середніх:

Оскільки метод k-середніх, як і багато інших методів кластеризації, має багато модифікацій, сконцентруємося перш за все на тому варіанті, що реалізований в цій роботі.

1.1) Постановка задачі:

Розглядаємо n спостережень $\bar{X} = (x_1, \dots, x_n)$, кожне з яких характеризується m ознаками. Серед інших алгоритмів кластеризації метод k-середніх виділяється тим, що ми наперед задаємо кількість кластерів k , на які буде розбита множина спостережень.

1.2) Алгоритм:

1.2.1) Задаємо k центроїдів – точок, що будуть тимчасовими центрами кластерів. Можна їх обрати випадковим чином або взяти конкретні елементи з множини спостережень, але зазвичай цей вибір робиться свідомо з урахуванням специфіки конкретної задачі, оскільки початковий вибір центроїдів має істотний вплив на роботу методу.

1.2.2) Перебираємо усі точки та знаходимо відстані до кожного з центроїдів. Для цього можна використовувати багато різних метрик, але на першому етапі використовується евклідова метрика $\rho(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$. Точка приєднується до того кластера чий центроїд є найближчим. Ця ітерація завершується на тому етапі, коли усі точки розподілені по кластерам.

1.2.3) Обираємо нові центроїди. Це також можна робити різними засобами, але в цій реалізації новий центроїд будемо обчислювати наступним чином:

$$c_i = \frac{1}{m_i} \cdot \sum_{x \in C_i} x, \text{ де } C_i \text{ – кластер; } c_i \text{ – його новий центроїд; } m_i \text{ – кількість точок,}$$

що були приєднані до цього кластера. Пізніше буде відзначено, що такий вибір засобу перерахунку центроїдів є оптимальним для евклідової метрики.

1.2.4) Виконуємо ті ж самі дії, що і на кроці 2, для нового набору центроїдів. Алгоритм зупиняється коли нове розбиття на кластери не відрізняється від того, що було отримано на попередній ітерації.

1.3) Обчислювальна складність:

1.3.1) Часова складність:

Найбільш корисною на практиці є оцінка, коли ми наперед обмежуємо кількість ітерацій алгоритму. Цей підхід справджується, оскільки необхідна для збіжності кількість ітерацій зазвичай невелика (щонайменше в порівнянні з об'ємом даних). Якщо вважати за елементарні операції арифметичні дії, то часова складність методу k-середніх – це $O(I \cdot k \cdot t \cdot n)$, де I – наперед задана кількість ітерацій, k – кількість кластерів, n – кількість спостережень, t – кількість ознак кожного спостереження [35]. В рамках цієї задачі доцільно розглядати обчислення відстані між елементами як елементарну операцію, і тоді часова складність дорівнює $O(I \cdot k \cdot n)$. Бачимо, що ця складність є лінійною відносно кількості спостережень, що разом з малою кількістю ітерацій для збіжності та простотою реалізації робить метод k-середніх швидким та гнучким для великого об'єму даних у порівнянні з багатьма іншими алгоритмами (навіть у тих випадках, коли збільшення об'єму даних зменшує ефективність цього методу, цю проблему можна вирішити за допомогою різноманітних модифікацій).

1.3.2) Просторова складність:

Метод k-середніх є також досить економним з точки зору пам'яті, бо окрім початкових даних для кластеризації зберігається лише k центроїдів, тому просторова складність буде $O((n + k) \cdot m)$, де k – кількість кластерів, n – кількість спостережень, m – кількість ознак кожного спостереження. [35]

1.4) Збіжність відносно цільової функції

Важливою частиною методу кластеризації є вибір цільової функції, що допомагає оцінити оптимальність деякого розбиття на кластери. Для методу k-середніх існує декілька комбінацій метрики, правила перерахування центроїдів та цільової функції, що гарантують збіжність відносно цієї цільової функції. Однією з таких комбінацій є евклідова метрика $\rho(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$, правило перерахування центроїдів $c_i = \frac{1}{m_i} \cdot \sum_{x \in C_i} x$ та цільова функція SSE (sum of the squared error) [35], яка для метрики ρ записується наступним чином:

$$SSE = \sum_{i=1}^k \sum_{x \in C_i} \rho(x, c_i)^2$$

Треба зауважити, що навіть в цьому випадку метод k-середніх гарантує збіжність лише до локального мінімуму цільової функції, оскільки невдалий вибір початкових центроїдів може значно погіршити якість кластеризації. [35, 36]

1.5) Переваги та недоліки:

Підіб'ємо підсумки та розглянемо головні позитивні та негативні особливості методу k-середніх.

Переваги:

- Гнучкість та простота реалізації [36]
- Швидка збіжність [36]

- Мала просторова складність
- Поширеність та велика кількість модифікацій

Недоліки:

- Необхідність наперед задавати кількість кластерів та початкові центроїди
- Ефективність методу дуже сильно залежить від вдалості вибору початкових центроїдів [36]
- Можна гарантувати збіжність лише до локального мінімуму цільової функції [35, 36]

Метод k-середніх був реалізований за допомогою програми в Python (вона наведена у додатку 3).

2) Метод найближчого сусіда

Оскільки метод найближчого сусіда не є основним в цій роботі, ми розглянемо його лише на базовому рівні. Перш за все, нам треба розуміти його алгоритм та мати змогу порівняти його з методом k-середніх.

2.1) Постановка задачі:

Розглядаємо n спостережень $\bar{X} = (x_1, \dots, x_n)$, кожне з яких характеризується m ознаками. На відміну від методу k-середніх, кількість кластерів не є наперед заданою. До того ж метод сусіду є ієрархічним агломеративним методом кластеризації (тобто ми починаємо з того, що кожне спостереження є окремим кластером, і потім крок за кроком з'єднуємо їх у кластери, просуваючись вгору за ієрархією), тому як результат ми маємо отримати ієрархічну структуру, яку зручно ілюструвати за допомогою дендрограми.

2.2) Алгоритм:

На кожному кроці обираються два найближчих спостереження (відстань між спостереженнями визначається за допомогою метрики, в даній реалізації це евклідова метрика $p(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$) і «зливаються» в одне (два часових ряди замінюються на один ряд, елементи якого є середніми арифметичним відповідних елементів цих 2 рядів). Кожне таке злиття фіксується задля подальшої побудови дендрограми. Повторюємо ці дії доки не залишається один часовий ряд. На основі послідовності «злиття» часових рядів та відстаней між ними будуюмо дендрограму, що наочно демонструє міру подібності між елементами.

2.3) Обчислювальна складність:

2.3.1) Часова складність:

В цьому дослідженні використана так званий «наївний алгоритм» для методу найближчого сусіда, що має часову складність $O(n^3)$, де n – кількість спостережень [37]. Існує декілька модифікацій, що можуть пришвидшити цей алгоритм, але найменша можлива часова складність для методу найближчого сусіда дорівнює $O(n^2)$ [38]. Можемо побачити, що при зростанні n метод найближчого сусіда буде працювати повільніше, ніж метод k -середніх, тому він гірше підходить для великих даних.

2.3.2) Просторова складність:

Оскільки в цій конкретній реалізації ми на кожному кроці будуюмо матрицю відстаней між кластерами, просторова складність буде також більшою ніж у методу k -середніх.

Метод найближчого сусіда був реалізований за допомогою програми у Python (вона наведена у додатку 3).

Незважаючи на такі переваги евклідової метрики як простота та гарантована збіжність за цільової функцією SSE (sum of the squared error), вона також має один істотний недолік, що може негативно впливати на коректність кластеризації: вона не враховує взаємні кореляції між величинами у багатомірному векторі (в нашій задачі ми іноді будемо розглядати часові ряди як вектори з деякого багатовимірного розподілу). Як вже було відзначено, цю проблему можна подолати за допомогою переходу від евклідової метрики до відстані Махаланобіса. Тому на другому етапі ми будемо спиратися на відстань Махаланобіса при застосуванні методу k-середніх, щоб провести порівняння з результатами, які були отримані на першому етапі, і зробити більш точні висновки щодо закономірностей, які спостерігаються в даних про захворюваність на COVID-19.

Відстань Махаланобіса обчислюється за наступною формулою: $d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^T S^{-1} (\vec{x} - \vec{y})}$, де S – це коваріаційна матриця багатовимірного розподілу, до якого належать вектори \vec{x} і \vec{y} . Звідси вже можна побачити, що застосування відстані Махаланобіса до двох векторів спирається на припущення, що ці вектори взяті з одного розподілу ймовірностей [39]. В рамках задачі про дослідження статистики COVID-19 у різних країнах це припущення можна сформулювати так, що пандемія відбувається в усіх країнах за деяким єдиним спільним законом. [41]

Перед застосуванням відстані Махаланобіса слід розглянути деякі особливості нашої задачі. По-перше, було б недоцільно обчислювати відстань Махаланобіса безпосередньо часових рядів, тому що в цьому випадку нам довелося б обчислити коваріаційну матрицю $N \times N$, а через те, що число N у нашій задачі дуже велике, обчислювальна складність також стає надто великою. Щоб уникнути цього, замість часових рядів будуть розглядатися їх наступні статистичні показники: середнє арифметичне часового ряду, мінімальне та максимальне значення в часовому ряді та коефіцієнт асиметрії. Коефіцієнт

асиметрії являє собою міру асиметрії ймовірнісного розподілу дійсної випадкової величини відносно її середнього значення та обчислюється за формулою $\Gamma_1 = \mu_3 / \sigma^3$, де μ_3 — це центральний момент третього порядку, а σ — це середнє квадратичне відхилення. [41]

Зв'язок між значенням коефіцієнта асиметрії та графіком функції щільності демонструється на наступному графіку:

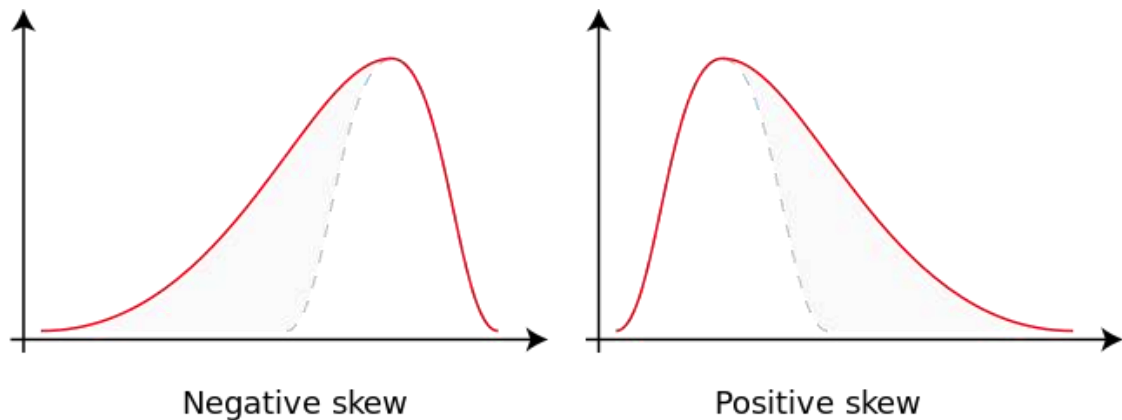


Рис. 5 Графік функції щільності в залежності від значення коефіцієнта асиметрії

Задля урівноваження впливу різних показників на значення відстані, деякі з них були пронормовані. Фінальний набір даних, до якого буде застосований метод k -середніх з відстанню Махаланобіса, виглядає наступним чином:

	Europe	Albania	Andorra	Austria	Belarus	Belgium	Bosnia	Bulgaria	Croatia	Czech	Denmark	Estonia	Finland	France
срзнач	0,4181124	0,3114335	1	0,4309021	0,3038822	0,5556092	0,3754154	0,392133	0,5270824	0,8042901	0,3218038	0,6486527	0,1349642	0,5780258
макс	0,2690729	0,2739007	0,9321128	0,5645353	0,1569713	1	0,3536472	0,3835695	0,6267098	0,8264677	0,4185573	0,8021729	0,0990633	0,5572296
мін	3,9165059	0,3414478	0,0016626	0,9778906	0,1442364	2,3727871	0,5412926	0,4460715	0,0192519	1,4090995	1,0043828	0,2079108	0,2981796	2,063287
скос	0,2674092	0,6618773	0,8174951	1,3627679	0,0262101	2,5409319	0,9422859	0,8769077	1,1504929	0,8792718	1,6419017	1,1854094	0,4861813	0,738863

Рис. 6 Новий набір даних, до якого буде застосований метод k -середніх з відстанню Махаланобіса

Вплив переходу до іншого набору даних на результати кластеризації буде детально проаналізовано в подальшому.

Оскільки в нашій задачі коваріаційна матриця загального ймовірнісного розподілу S невідома, замість неї була використана вибіркова коваріаційна матриця (в якості вибірки був використаний весь новий набір даних). Вибіркова коваріаційна матриця $S = [s_{jk}]$ була обчислена за допомогою функції `numpy.cov()` за наступною формулою:

$s_{jk} = \frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$, де N – це розмір відповідної вибірки, \bar{x}_j та \bar{x}_k – відповідні вибіркові середні

Фінальним кроком в цій роботі є застосування методу k -середніх з евклідовою метрикою та з відстанню Махаланобіса до кумулятивних даних з джерела [32] (тобто даних про загальну кількість хворих на COVID-19 в країнах Європи). Загальна ідея кластерного аналізу та ж, що і на попередніх етапах: спочатку застосовується метод k -середніх з евклідовою метрикою, потім створюється новий набір даних з чотирма статистичними показниками і вже до нього застосовується метод k -середніх з відстанню Махаланобіса, наприкінці порівнюються результати цих двох засобів кластеризації.

5. Результати

На першому етапі після застосування методу k -середніх з початковими центроїдами, що були отримані за допомогою кластерів з дослідження [18] (обчислюємо центроїди як середні арифметичні часових рядів для країн, що входять до нього, як і було в описі методу. У другому кластері не беремо піддериво з Ліхтенштейну, Швейцарії, Люксембургу та Хорватії, оскільки до

нього входять дуже малі країни, які мають дуже специфічну статистику захворюваності і тому можуть негативно вплинути на роботу методу), і обчислення радіусів кластерів та відстаней до центру, отримуємо наступні 4 кластери (Таблиця 1).

Одразу можемо порівняти склад отриманих кластерів з тими, що були в дослідженні [18]. Для цього стане в нагоді наведена нижче порівняльна діаграма для результатів кластеризації методом найближчого сусіда та методом k-середніх (Рис.7).

Таблиця 1. Чотири кластери, що були отримані методом k-середніх з початковими центроїдами на основі результатів дослідження [18], відстані до центроїдів та радіуси кластерів

Кластер 1			Кластер 2		
Країни	Відстань до центру	Радіус	Країни	Відстань до центру	Радіус
Монако	2019,54	2933,8	Швейцарія	2418,38	5399,47
Іспанія	2092,76		Нідерланди	2544,36	
Ірландія	2452,35		Словенія	2738,23	
Португалія	2813,49		Люксембург	2742,08	
Великобританія	2933,8		Франція	3094,83	
			Ліхтенштейн	3144,83	
			Бельгія	3474,64	
			Литва	4746,46	
			Андорра	5081,48	

Кластер 3			Кластер 4		
Країни	Відстань до центру	Радіус	Країни	Відстань до центру	Радіус
Болгарія	1282,05	5373,97	Білорусь	816,13	1850,46
Македонія	1631,92		Німеччина	1027,87	
Боснія	1773,06		Росія	1053,73	
Австрія	1808,17		Норвегія	1195,49	
Італія	1882,26		Фінляндія	1475,34	
Молдова	1898,58		Україна	1599,45	
Польща	2378,5		Греція	1746,58	
Румунія	2404,54		Данія	1802,82	
Словаччина	2629,36		Албанія	1850,46	
Угорщина	2682,12				
Хорватія	2973,06				
Швеція	2976,58				
Латвія	3307,32				
Косово	4160,92				
Естонія	4551,05				
Сербія	5373,97				
			Чехія	5399,47	

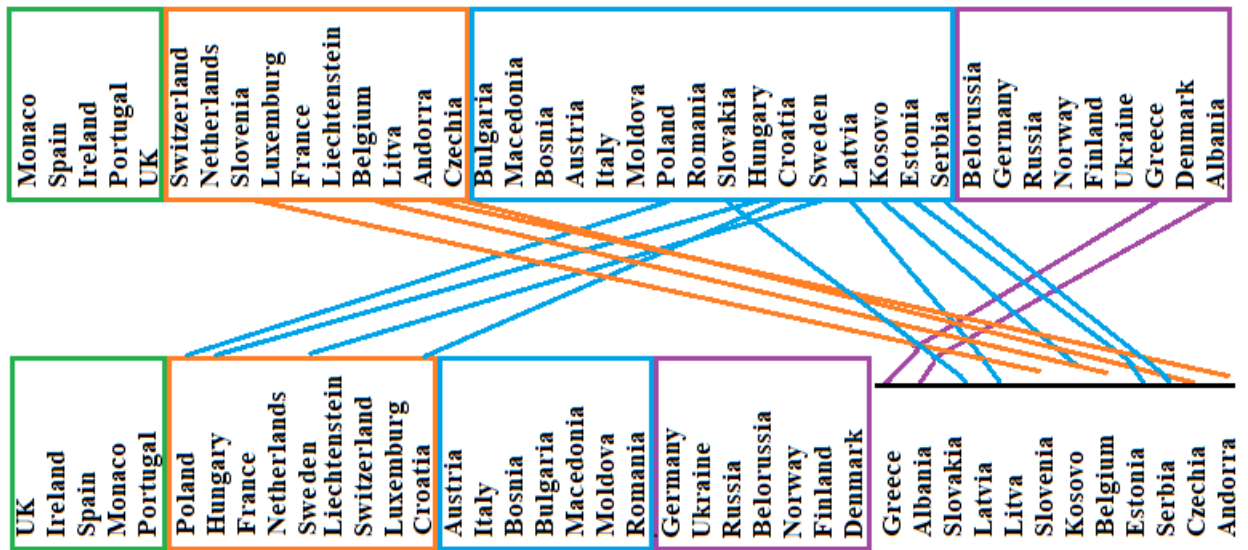


Рис. 7 Співставлення кластерів, що були знайдені в дослідженні [18] за допомогою методу найближчого сусіда (внизу) та кластерів, що були знайдені методом k-середніх (лінії позначають перехід з одного кластера в інший)

Подивимось перш за все на країни, на основі яких були сформовані центроїди. Тут результати досить добре збігаються з методом найближчого сусіду: лише Польща, Угорщина та Швеція перейшли з 2 до 3 кластеру, інші ж країни залишилися там, де і були. Цікаво, що Ліхтенштейн, Люксембург та Швейцарія також залишилися в 2 кластері, хоча вони і не використовувалися для обчислення центроїдів. Інші країни стояли на рис. 3 як би окремо від чотирьох кластерів, але при застосуванні методу k-середніх вони розподілились між цими кластерами.

Тепер задля того, щоб представити отримані кластери у вигляді дерев та порівняти їх внутрішню структуру з кластерами з дослідження [18], застосовуємо метод найближчого сусіда до кожного з кластерів, тільки в цьому випадку оберемо в якості метрики модуль різниці між відстанями до центру відповідного кластера. Отримуємо наступні дендрограми (номер біля дендрограми – номер

відповідного кластера) (Рис.8). Бачимо, що структура цих кластерів дуже відрізняється від тих, що були в дослідженні [18].

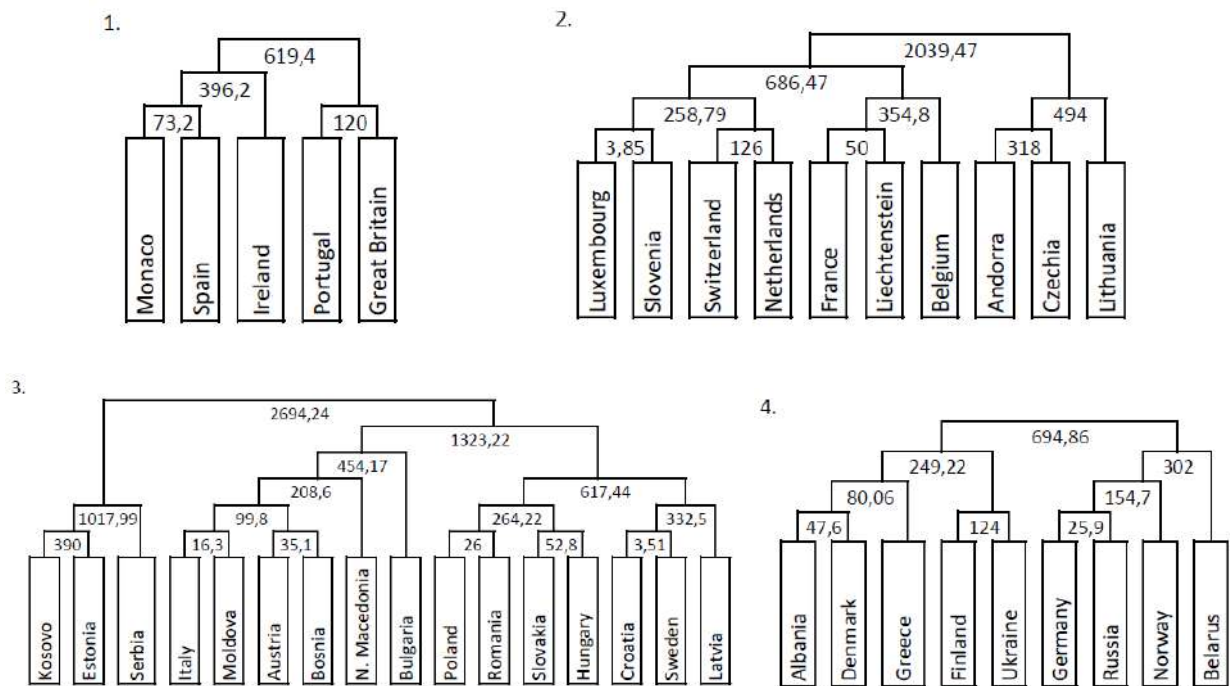


Рис. 8 Дендрограми для 4 кластерів, що були отримані за допомогою методу k-середніх (початкові центроїди на основі результатів дослідження [18])

Ми вже відмічали, що вибір початкових центроїдів має дуже великий вплив на результати роботи методу k-середніх, тому спробуємо побудувати кластеризацію спираючись на інші початкові центроїди. Для цього розіб'ємо усю вибірку європейських країн (за винятком дуже малих країн, оскільки вони мають дуже специфічну статистику і це може негативно вплинути на роботу методу k-середніх) на 7 груп країн, що близькі одна до одної, та обчислимо центроїди як середні арифметичні часових рядів для країн, що входять до відповідної групи, як і було в описі методу k-середніх. [41]

Групи географічно близьких країн:

Група 1: Іспанія, Португалія

Група 2: Франція, Бельгія, Нідерланди

Група 3: Великобританія, Ірландія

Група 4: Німеччина, Австрія, Швейцарія, Італія

Група 5: Данія, Швеція, Норвегія, Фінляндія

Група 6: Польща, Україна, Естонія, Латвія, Литва, Чехія, Словаччина, Угорщина

Група 8: Словенія, Хорватія, Боснія, Сербія, Македонія, Албанія, Греція, Румунія, Молдова

Отримуємо наступні 7 кластерів:

Таблиця 2. Чотири кластери, що були отримані методом k-середніх з початковими центроїдами на основі географічних груп, відстані до центроїдів та радіуси кластерів

Кластер 1			Кластер 2		
Країни	Відстань до центру	Радіус	Країни	Відстань до центру	Радіус
Іспанія	1436,7	2198,95	Франція	2412,36	4434,38
Монако	1892,58		Нідерланди	2553,26	
Португалія	2198,95		Бельгія	2628,78	
		Андорра	4434,38		
Кластер 3			Кластер 4		
Країни	Відстань до центру	Радіус	Країни	Відстань до центру	Радіус

Ірландія	1700,35	1700,35	Австрія	1212,4	2906,48
Великобританія	1700,35		Італія	1446,75	
			Швейцарія	1632,44	
			Боснія	2055,13	
			Польща	2245,91	
			Люксембург	2693,95	
			Ліхтенштейн	2906,48	
Кластер 5			Кластер 6		
Країни	Відстань до центру	Радіус	Країни	Відстань до центру	Радіус
Беларусь	861,58	1908,65	Словаччина	2445,08	5191,29
Німеччина	1045,53		Швеція	2613,18	
Росія	1139,71		Словенія	3213,89	
Норвегія	1311,43		Латвія	3311,61	
Україна	1464,43		Угорщина	3320,54	
Фінляндія	1586,58		Естонія	3767,42	
Албанія	1766,09		Чехія	5191,29	
Греція	1771,89				
Данія	1830,83				
Молдова	1908,65				
Кластер 7					

Країни	Відстань до центру	Радіус
Македонія	1976,87	4423,76
Болгарія	2016,37	
Хорватія	2394,35	
Румунія	2640,43	
Литва	3951,12	
Косово	4136,76	
Сербія	4423,76	

Одразу можемо порівняти склад кластерів зі складом початкових груп. Бачимо дуже багато відмінностей: хоча менші кластери (1, 2, 3,7) змінилися не так сильно, склад більших кластерів (4, 5, 6) вже має мало спільного з початковими групами.

Тепер, як і в попередньому випадку, застосуємо метод найближчого сусіда до отриманих кластерів із метрикою у вигляді модулю різниці відстаней до центру відповідних кластерів. В якості результату отримаємо дендрограму на рис. 12, на якій можна побачити, що внутрішня структура кластерів, побудованих за географічними групами, кардинально відрізняється від структури кластерів на рис. 3. Ця дендрограма та мапи, які допомагають відстежувати географічні закономірності в результатах кластеризації, наведені в додатку 1.

Треба зауважити, що кластер 3 має цікаву особливість: оскільки він складається з двох країн і його центроїд є середнім арифметичним часових рядів для цих країн, вони мають однакові відстані до центру. Тому якщо ми будемо

розраховувати відстань між елементами як модуль різниці відстаней до центра, отримаємо 0, що, звісно, не відповідає дійсності. Цей приклад показує недолік відстані, що була введена заради побудови дендрограм методом найближчого сусіда для кластерів, які були отримані за допомогою метода k-середніх.

На другому етапі застосовуємо метод k-середніх із відстанню Махаланобіса. Слід враховувати, що перехід до іншого набору даних може істотно вплинути на результати кластеризації. Тому для нового набору даних зі статистичними показниками використовуються як декартова метрика, так і відстань Махаланобіса, щоб відокремити вплив зміни метрики від впливу зміни набору даних. Як і на першому етапі, будемо кластеризацію за допомогою метода k-середніх з 4 початковими центроїдами на основі результату методу найближчого сусіда та з 7 початковими центроїдами на основі географічних груп європейських країн. Результати кластеризації зручно проілюструвати за допомогою кольорових мап (Рис.9-10, Cartesian metrics означає евклідову метрику).

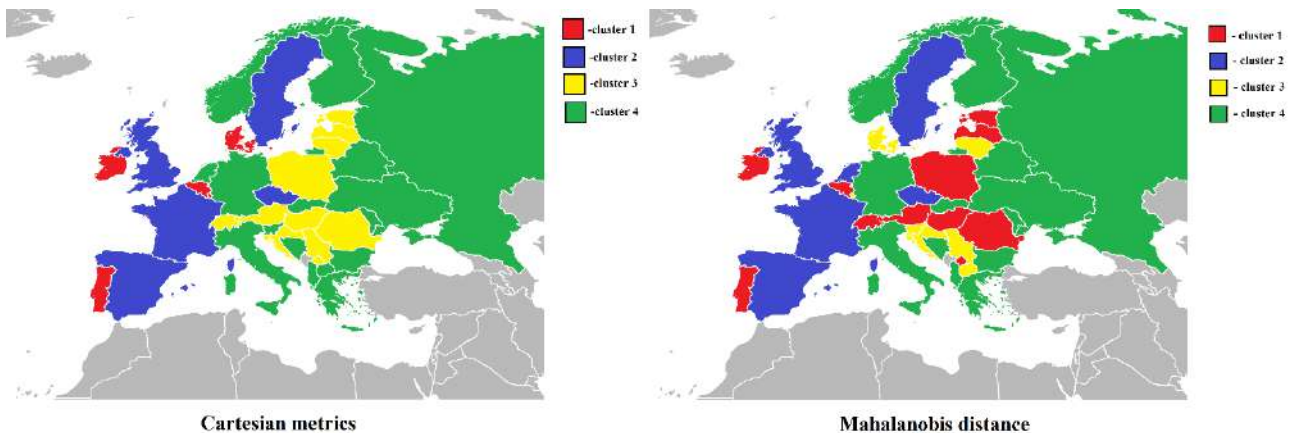


Рис. 9 Візуалізація на мапі Європи 4 кластерів (з початковими центроїдами на основі результатів метода найближчого сусіда) отриманих при застосуванні метода k-середніх до набору даних зі статистичними показниками

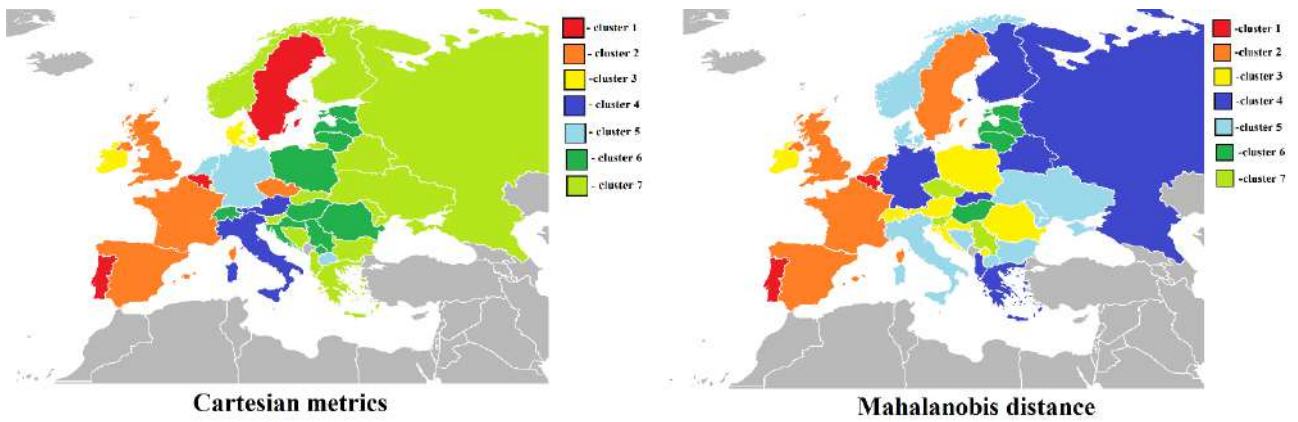


Рис. 10 Візуалізація на мапі Європи 7 кластерів (з початковими центроїдами на основі груп європейських країн) отриманих при застосуванні метода k -середніх до набору даних зі статистичними показниками

Після порівняння результатів кластеризації нового набору даних із попередніми результатами кластеризації ми можемо зробити такі висновки: в нових результатах стало менше відповідності між географічною близькістю країн та подібністю їхньої динаміки захворюваності на COVID-19, але більшість основних змін пояснюється переходом до іншого набору даних. Що стосується наслідків зміни метрики, то вона нечасто розбиває групи країн, що раніше були разом в одному кластері, натомість може помістити їх усіх разом в інший кластер. Наприклад, ми бачимо на рис. 9, що група з Польщі, Швейцарії, Австрії, Угорщини та Румунії належить до одного кластеру як після застосування декартової метрики, так і після застосування відстані Махаланобіса, але в першому випадку вона належить до кластера 3, а в останньому випадку він належить до кластера 1. [41]

Для оцінки впливу географічного фактору також корисно розглянути розмір найбільшої групи країн зі спільними кордонами в кластері для різних результатів кластеризації. Відповідні дані містяться в таблицях 3 та 4, які

наведені в додатку 2. Вони демонструють, що загалом перехід від набору даних с часовими рядами до набору даних з їх статистичними показниками та від евклідової метрики до відстані Махаланобіса призводить до зникнення географічних закономірностей, що були виявлені раніше, та робить результати кластеризації складнішими для інтерпретації.

Незважаючи на те, що результати кластеризації на другому етапі показують, що відповідність між географічною близькістю країн та їхньою схожістю за статистикою захворюваності на коронавірус є досить нестійкою, все ж є деякі закономірності, які представляють великий інтерес для подальших досліджень. Можна помітити декілька груп країн, які географічно не близькі та на перший погляд мають небагато спільних рис, але дуже часто потрапляють до одного кластера, наприклад Франція та Чехія (вони належать до одного кластера в усіх випадках, окрім результатів кластеризації з 7 географічними центроїдами з використанням евклідової метрики для часових рядів та відстані Махаланобіса для набору даних зі статистичними показниками) або Німеччина, Білорусь та Росії (вони належать до одного кластеру в усіх випадках, окрім результатів кластеризації з 7 географічними центроїдами з використанням евклідової метрики набору даних зі статистичними показниками).

В цій ситуації варто було б окремо перевірити, чи дійсно країни в цих нетипових групах схожі за статистикою COVID-19. Щоб оцінити їхню подібність, ми будемо два окремих графіки для цих груп (на першому графіку ми також додаємо криву для Нідерландів, оскільки ця країна часто знаходиться в одному кластері з Францією та Чехією, але географічно близька до Франції і тому інтуїтивно здається більш схожою на неї) і один спільний графік для обох груп.

На цих графіках можна побачити, що криві для країн з однієї групи досить схожі за піками захворюваності та їх крутизною (піки захворюваності Білорусі, Німеччини та Росії зміщені вперед порівняно з піками захворюваності в Чехії,

Франції та Нідерландах і є значно менш крутими), тоді як криві для країн з різних груп помітно відрізняються. Це спостереження доводить достовірність отриманих результатів кластеризації. Усі три порівняльні графіки наведені в додатку 1.

На останньому етапі застосовуємо метод k-середніх із евклідовою метрикою та відстанню Махаланобіса до кумулятивних даних (тобто часових рядів із загальною кількістю нових випадків на мільйон). Як і на другому етапі, для відстані Махаланобіса використовується набір даних із чотирма статистичними показниками. Якщо знов проілюструвати результати кластеризації за допомогою кольорових мап, то можна побачити, що вони демонструють не більш стійкий зв'язок між географічною близькістю та статистикою COVID-19, ніж на попередніх етапах. Однак на усіх чотирьох мапах чітко виділяється новий кластер зі східноєвропейських і балканських країн (він зображений окремо на наступній мапі). Усі 4 мапи з результатами кластеризації кумулятивних даних для різних початкових центроїдів та метрик можна знайти в додатку 1.



Рис. 11 Кластер зі східноєвропейських і балканських країн, що спостерігається в результатах кластеризації кумулятивних даних

Незважаючи на те, що деякі з цих країн часто потрапляли в один кластер в попередніх результатах кластеризації, не було об'єднаного географічного блоку країн, який вже можна чітко побачити чітко на мапах для кумулятивних даних. Оскільки цей кластер не розпадається після зміни початкових центроїдів і метрики, подібність між цими країнами можна вважати сильною та стійкою. Цей висновок підтверджується графіком загальної кількості випадків на мільйон у Боснії, Польщі та Румунії (ці три країни були взяті як вибірка з кластера) та графіком для усіх країн з кластера, які наведені у додатку 1. Помітно, що ці три країни мають багато спільного у динаміці захворюваності, а графіки в цілому дуже близькі один до одного, особливо в часовому проміжку з березня 2020 року до березня 2021 року.

6. Висновки

Результати кластеризації методом k -середніх в цілому збігаються з результатами, отриманими методом найближчого сусіда (принаймні для тих країн, які були враховані при розрахунку початкових центроїдів). Оскільки результати, отримані методом найближчого сусіда, вважаються досить надійними та коректними, можна стверджувати про коректність результатів метода k -середніх із зазначеним вибором початкових центроїдів. Але структура дендрограм, які були побудовані за результатами кластеризації за допомогою метода k -середніх, вже суттєво відрізняється від структури дендрограм, які були побудовані за методом найближчого сусіда. Крім того, обрана для побудови нових дендрограм відстань не завжди дає правильні результати, як це було з двоелементним кластером. Отже, можна зробити висновок, що для визначення

складу кластерів доцільно використовувати метода k-середніх, але для побудови ієрархічних структур краще окремо використовувати інші методи. [41]

Щоб зробити результати кластеризації більш правильними завдяки урахуванню кореляцій у наборі даних, був застосований метод k-середніх з метрикою Махаланобіса. Задля того, щоб цей метод можна було раціонально застосувати, довелося перейти до нового набору даних зі статистичними показниками, і виявилось, що цей перехід набагато більше впливає на результати кластеризації, ніж зміна метрики. [41]

Великий інтерес представляє питання того, наскільки схожими є географічно близькі країни за статистикою захворюваності. Аналізуючи результати метода k-середніх, можна сказати, що географічно близькі країни часто потрапляють в один кластер, навіть частіше, ніж при використанні метода найближчого сусіда. Тим не менш, виявлені географічні закономірності часто є нестабільними і можуть зникнути при переході до нового набору даних або при зміні евклідової метрики на відстань Махаланобіса. Також у результатах обох методів ми можемо побачити багато випадків, коли країни, які знаходяться далеко одна від одної, потрапляють в один кластер, і деякі з цих незвичайних сполук утворюються як при застосуванні метода k-середніх (і майже всіх його реалізованих варіацій), так і при застосуванні метода найближчого сусіда: наприклад, Росія, Німеччина і Білорусь часто потрапляють в один кластер (зазвичай разом з кількома країнами Північної Європи). Такі незвичайні сполуки навряд чи є цілком випадковими і тому потребують подальшого дослідження. Однак кластерний аналіз кумулятивних даних (тобто часових рядів для загальних випадків захворюваності на мільйон населення) виділяє стабільний і географічно однорідний кластер східноєвропейських і балканських країн, який відповідає географічній логіці набагато краще, ніж попередні сполуки. Порівняння динаміки захворюваності на коронавірус у деяких країнах з цих сполук довело, що вони

справді схожі за статистикою захворюваності та що результати кластеризації є коректними. [41]

Отримані результати показують, що географічна близькість країн має значний вплив на їх статистику захворюваності, але є й інші важливі фактори, які пояснюють подібність статистики коронавірусу в географічно віддалених країнах. [41]

7. Список використаних джерел:

1. Covid-19 Pandemic. Problems Arising in Health and Social Policy. /Ch. Aspalter (Ed.). Springer Singapore. 2023. 303 p. <https://doi.org/10.1007/978-981-99-2497-4>
2. Koley T.K., Dhole M. The COVID-19 Pandemic: The Deadly Coronavirus Outbreak. Taylor & Francis, 2020. - 174 p.
3. Nandiyanto A.B.D. COVID-19: From Health, Education, Economic, to Science and Technology in South East Asia and India. Nova Science Publ., Inc. 2021. - 385 p.
4. Cascella M., Rajnik M., Aleem A., et al. Features, Evaluation, and Treatment of Coronavirus (COVID-19). In: StatPearls [Internet]. Treasure Island: StatPearls Publishing. 2024. Available from: <https://www.ncbi.nlm.nih.gov/books/NBK554776/>
5. Spiegelhalter D., Masters A. Covid by numbers: Making sense of the pandemic with data. Penguin UK, 2022. - 320 p. <https://doi.org/10.1093/ocmed/kqac013>
6. Zakharova A.A., Kizilova N.M. Study of correlations of the dynamics of COVID-19 disease with some socio-economical factors. Bull. of V.N. Karazin Kharkiv National University. ser. "Mathematical modeling. Information Technology. Automated control systems". – 2020. – N48. – P.49-56. (in Ukrainian)
7. Fox W.P. Mathematical Modeling in the Age of the Pandemic. CRC Press, Taylor & Francis Group, 2021. 164 p.
8. Vyta V., Ramakuri S.K., Peddi A., et al Mathematical Models for Predicting Covid-19 Pandemic: A Review. J. Phys.: Conf. Ser. 2020. Vol.1797, 012009. doi: 10.1088/1742-6596/1797/1/012009
9. Krishna M.V. Mathematical modelling on diffusion and control of COVID–19. Infectious Disease Modelling. 2020, Vol.5, P.588-597. <https://doi.org/10.1016/j.idm.2020.08.009>
10. Kizilova N. Nonlinear Dynamics, Stability and Control Strategies: Mathematical Modeling on the Big Data Analyses of Covid-19 in Poland. // 16th Intern. Conf. "Dynamical Systems: Theory and Applications". Lodz University of Technology. - 2021. - P.467-468.

11. Fraguera-Collar A. Moving From COVID-19 Mathematical Models to Vaccine Design: Theory, Practice and Experiences. Bentham Sci. Publ. 2022. 583 p.
12. Brauer, F., Mathematical epidemiology: Past, present, and future. *Infect. Dis. Model.* 2(2), 113–127 (2017).
13. Overton Ch.E., Stage H.B., Ahmad S., et al. Using statistics and mathematical modelling to understand infectious disease outbreaks: COVID-19 as an example. *Infectious Disease Modelling*, 2020, Vol.5. P.409-441. <https://doi.org/10.1016/j.idm.2020.06.008>.
14. He X., Lau E.H.Y., Wu P., et al. Temporal dynamics in viral shedding and transmissibility of COVID-19, *Nat. Med.* 2020, Vol. 26, P. 672–675.
15. Chun, J.Y., Baek, G., Kim, Y., Transmission onset distribution of COVID-19, *Intern. J. Infect. Dis.*, 99, 403–407 (2020).
16. Danelon A.F., Kumbhakar S.C. Estimating COVID-19 under-reporting through stochastic frontier analysis and official statistics: A case study of São Paulo State, Brazil. *Socio-Economic Planning Sciences*, 2023. Vol. 90, 101753. <https://doi.org/10.1016/j.seps.2023.101753>.
17. Kostecka V.V., Kizilova N.M. Mathematical modeling of dynamics of the covid-19 epidemic. *Bull. of V.N. Karazin Kharkiv National University. ser. "Mathematical modeling. Information Technology. Automated control systems"*. – 2020. – N48. – P.65-71. (in Ukrainian)
18. Kulik D.A., Kurkchi Ye.P., Kizilova N.M. Analysis of "big data" and mathematical modeling of the covid-19 epidemic in European countries. *Bull. of V.N. Karazin Kharkiv National University. ser. "Mathematical modeling. Information Technology. Automated control systems"*. – 2021. – N52. – P.35-42. (in Ukrainian) <https://doi.org/10.26565/2304-6201-2021-52-05>
19. Voloshyna K.I., Kizilova N.M., Kiporenko P.V. Study of the dynamics of four waves of covid-19 in European countries. *Bull. of V.N. Karazin Kharkiv National*

University. ser. "Mathematical modeling. Information Technology. Automated control systems". – 2021. – N54. – P.6-15. (in Ukrainian)

20. Sempé L., P. Lloyd-Sherlock, R. Martínez, S. Ebrahim, M. McKee, E. Acosta Estimation of all-cause excess mortality by age-specific mortality patterns for countries with incomplete vital statistics: a population-based study of the case of Peru during the first wave of the COVID-19 pandemic. *The Lancet Regional Health - Americas*, 2021, Vol.2, 100039. <https://doi.org/10.1016/j.lana.2021.100039>.

21. Aldila, S.H.A., Khoshnaw, E., Safitri, et al., A mathematical study on the spread of COVID-19 considering social distancing and rapid assessment: The case of Jakarta, Indonesia, *Chaos, Solitons & Fractals*, 139, 110042 (2020).

22. Li Q., Guan X., Wu P., et al. Early transmission dynamics in Wuhan, China, of novel coronavirus–infected pneumonia. *New England Journal of Medicine*, 2020. Vol. 382, 1199-1207. <https://doi.org/10.1056/NEJMoa2001316>

23. Asadi-Zeydabadi M., Buscema M., Lodwick W., et al. Analysis of COVID-19 pandemic in USA, using Topological Weighted Centroid. *Computers in Biology and Medicine*. 2021. Vol.136, 104670. <https://doi.org/10.1016/j.compbimed.2021.104670>

24. Sun J. Forecasting COVID-19 pandemic in Alberta, Canada using modified ARIMA model. *Computer Methods and Programs in Biomedicine*. 2021. Vol. 22, 100029. <https://doi.org/10.1016/j.cmpbup.2021.100029>

25. Galanis G., Georgiadis A. Socioeconomic conditions and contagion dynamics of the COVID-19 pandemic with and without mitigation measures: Evidence from 185 countries. *World Development*, 2024, Vol.175, 106477. <https://doi.org/10.1016/j.worlddev.2023.106477>.

26. Žid P., M. Haindl, V. Havlíček Governmental Anti-Covid Measures Effectiveness Detection. *Procedia Computer Science*, 2023, Vol.225, P.2922-2931. <https://doi.org/10.1016/j.procs.2023.10.285>.

27. Tomer V., Gupta S., Manwal M., Singh D.P. How statistics of World health index react against COVID-19. *Materials Today: Proceedings*, 2021. Vol.46, P.11267-11273. <https://doi.org/10.1016/j.matpr.2021.03.486>.
28. de Araújo Morais L.R., da Silva Gomes G.S. Applying Spatio-temporal Scan Statistics and Spatial Autocorrelation Statistics to identify Covid-19 clusters in the world - A Vaccination Strategy? *Spatial and Spatio-temporal Epidemiology*. 2021. Vol.39. 100461. <https://doi.org/10.1016/j.sste.2021.100461>.
29. Ramadan A., Kamel A., Taha A., et al. A multivariate data analysis approach for investigating daily statistics of countries affected with COVID-19 pandemic. *Heliyon*, 2020. Vol.6, N11, e05575. <https://doi.org/10.1016/j.heliyon.2020.e05575>.
30. Na J., Tibebu H., De Silva V., Kondoz A., Caine M. Probabilistic approximation of effective reproduction number of COVID-19 using daily death statistics. *Chaos, Solitons & Fractals*, 2020. Vol.140,110181. <https://doi.org/10.1016/j.chaos.2020.110181>.
31. *Big Data Analysis: New Algorithms for a New Society*, ed. by N. Japkowicz, J. Stefanowski, Springer (2016).
32. Research and data to make progress against the world's largest problems. [Electronic resource]. Available: <https://ourworldindata.org/coronavirus>
33. Ahmed Hamed, Ahmed Sobhy, Hamed Nassar. Accurate Classification of COVID-19 Based on Incomplete Heterogeneous Data using a KNN Variant algorithm. *Arabian Journal for Science and Engineering* 2021. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7931985/>
34. Кізілова Н.М. Методи аналізу «великих даних»: методичні рекомендації з курсу «Прикладні задачі аналізу великих даних» для здобувачів вищої освіти спеціальності «Прикладна математика». – Харків : Харківський національний університет імені В.Н. Каразіна, 2023.

35. Pan-Ning Tan, Michael Steinbach, Anuj Karpatne, Vipin Kumar. Introduction to Data Mining (Second Edition). 2019. Chapter 7. P. 10-20
36. О. М. Ткаченко, О. Ф. Грійо Тукало, О. В. Дзісь, С. М. Лаховець. Метод кластеризації на основі послідовного запуску k-середніх з удосконаленим вибором кандидата на нову позицію вставки. *НаукПраці ВНТУ*. 2012. С. 1
37. Murtagh F, Contreras P (2012). "Algorithms for hierarchical clustering: an overview". *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*. 2 (1). Wiley Online Library 86–97. doi:10.1002/widm.53
38. Gower JC, Ross GJ (1969). "Minimum spanning trees and single linkage cluster analysis". *Journal of the Royal Statistical Society, Series C*. 18 (1): 54–64. doi:10.2307/2346439. JSTOR 2346439. MR 0242315.
39. De Maesschalck, R.; Jouan-Rimbaud, D.; Massart, D. L. (2000). "The Mahalanobis distance". *Chemometrics and Intelligent Laboratory Systems*. 50 (1): 1–18. doi:10.1016/s0169-7439(99)00047-7
40. Фальченко І. Р. Кізілова Н. М. Аналіз «великих даних» динаміки COVID-19 у країнах Європи методами кластеризації з різними метриками. Збірник тез доповідей XVIII Міжнародної науково-практичної конференції студентів та молодих вчених «Сучасні проблеми математики та її застосування в природничих науках та інформаційних технологіях». – Харків : Харківський національний університет імені В.Н. Каразіна, 10-11 травня 2024р. С. 45-48
41. Falchenko I. R., Kizilova N. M. Cluster Classification with Different Metrics for Big Data Analysis of Covid-19 Dynamics. *Proceedings ICTERI 2024* (submitted)

8. Додаток 1

У цьому додатку наведені дендрограми, кольорові мапи з кластерами та графіки, які були згадані у цій роботі.



Рис. 12 Дендрограми для 7 кластерів, що були отримані за допомогою методу k -середніх (початкові центроїди на основі груп географічно близьких країн)

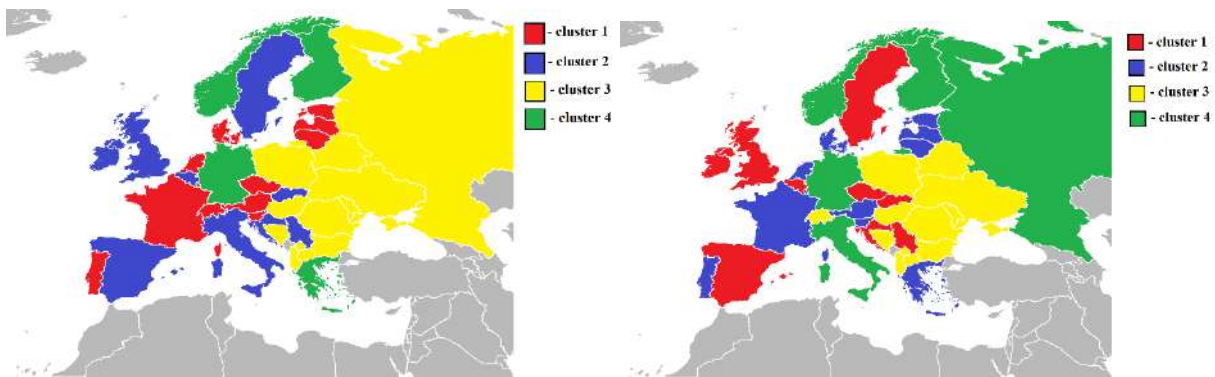


Рис. 13 Візуалізація на мапі Європи 4 кластерів (з початковими центроїдами на основі результатів метода найближчого сусіда) отриманих при застосуванні метода k -середніх до набору даних зі статистичними показниками (кумулятивні дані)

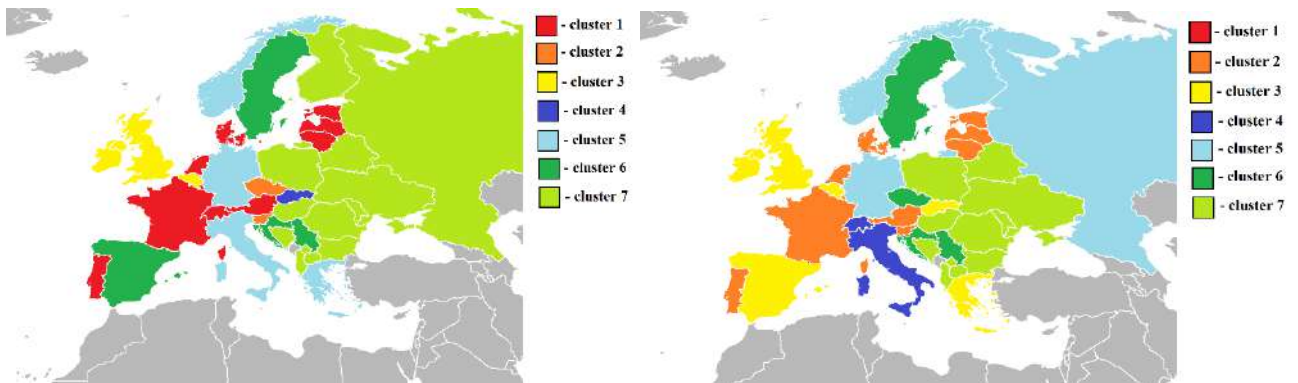


Рис. 14 Візуалізація на мапі Європи 7 кластерів (з початковими центроїдами на основі груп європейських країн) отриманих при застосуванні метода k -середніх до набору даних зі статистичними показниками (кумулятивні дані)

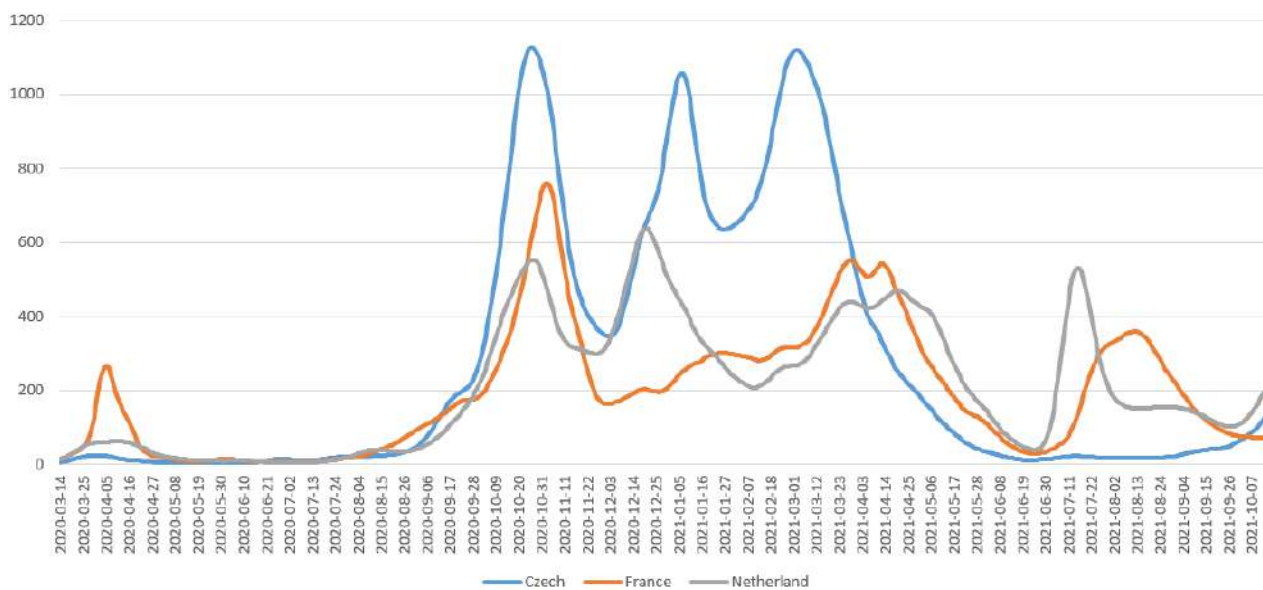


Рис. 15 Порівняльний графік нових випадків захворювання на мільйон населення для Франції, Чехії та Нідерландів

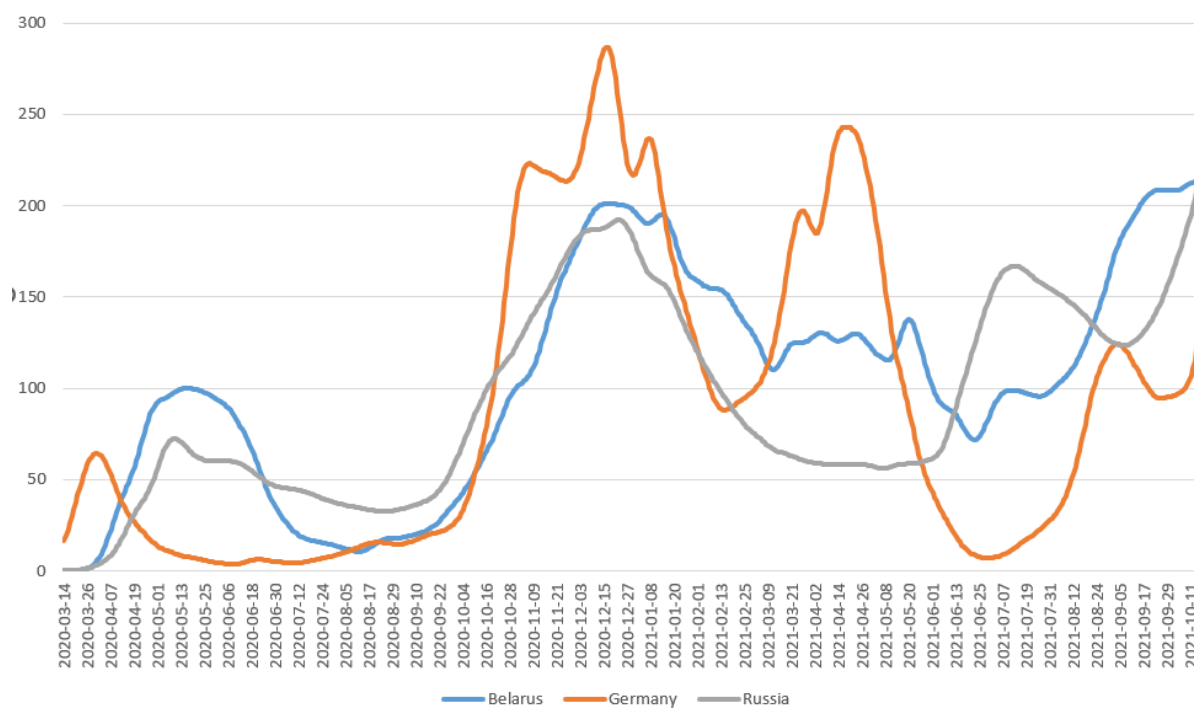


Рис. 16 Порівняльний графік нових випадків захворювання на мільйон населення для Беларусі, Німеччини та Росії

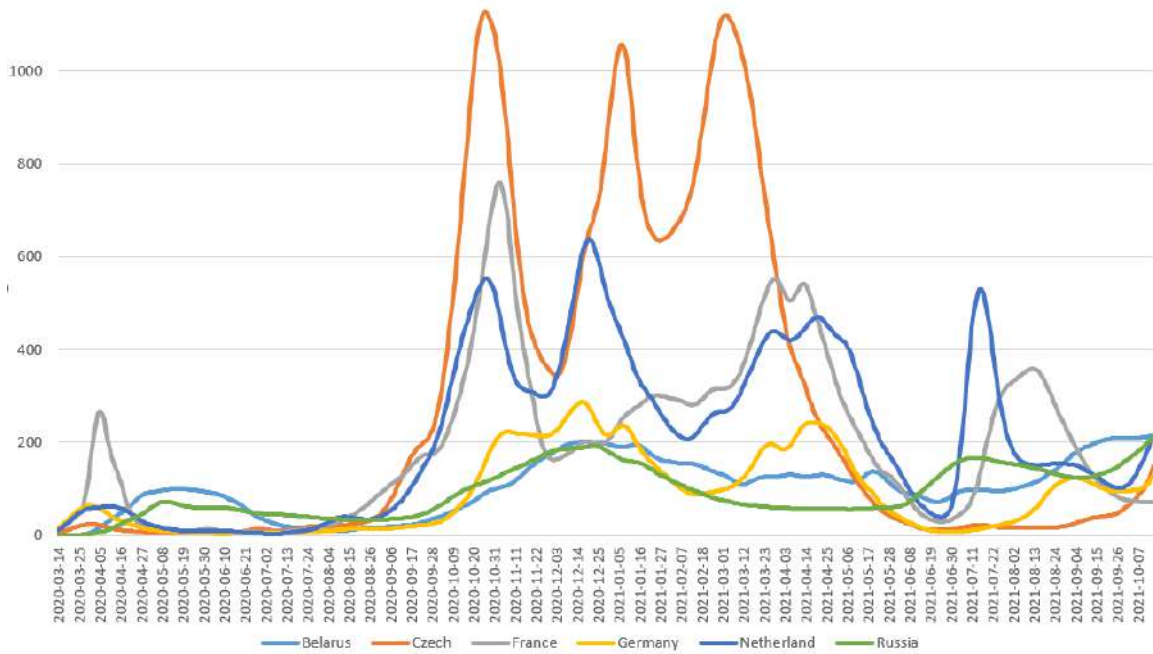


Рис. 17 Загальний порівняльний графік нових випадків захворювання на мільйон населення для Беларусі, Німеччини, Росії, Франції, Чехії та Нідерландів

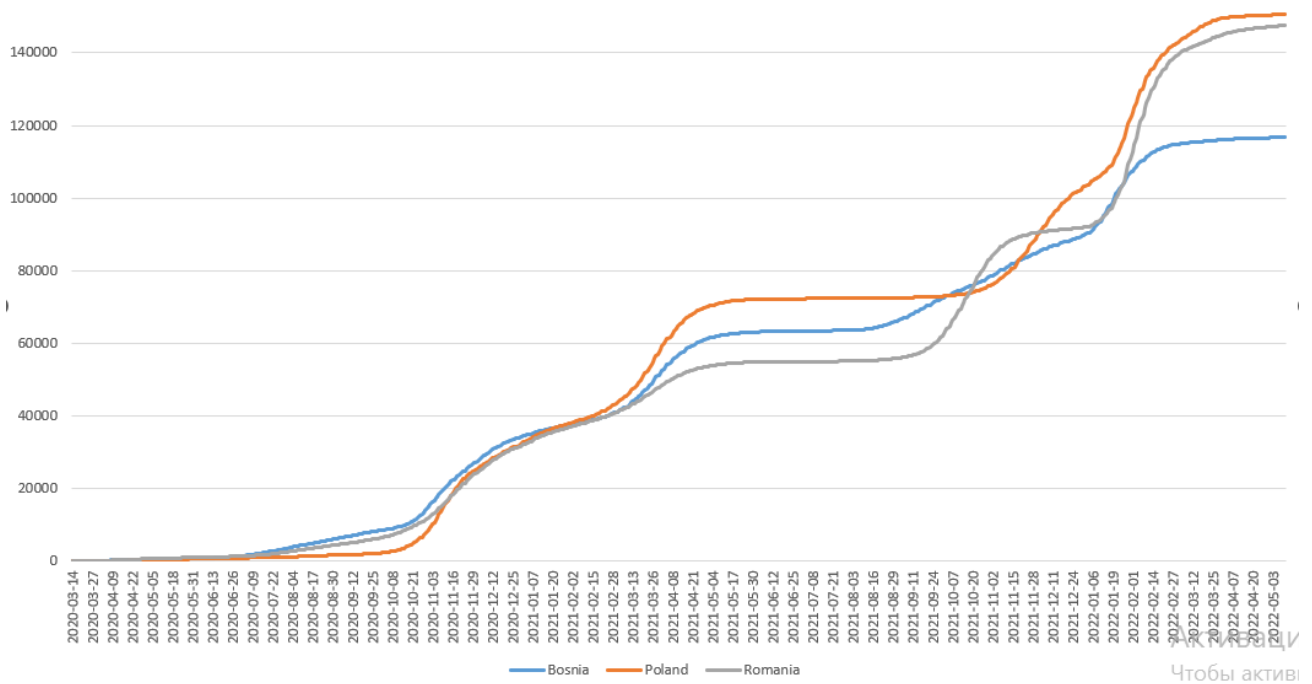


Рис. 18 Порівняльний графік загальної кількості випадків захворювання на мільйон населення для Боснії, Польщі та Румунії

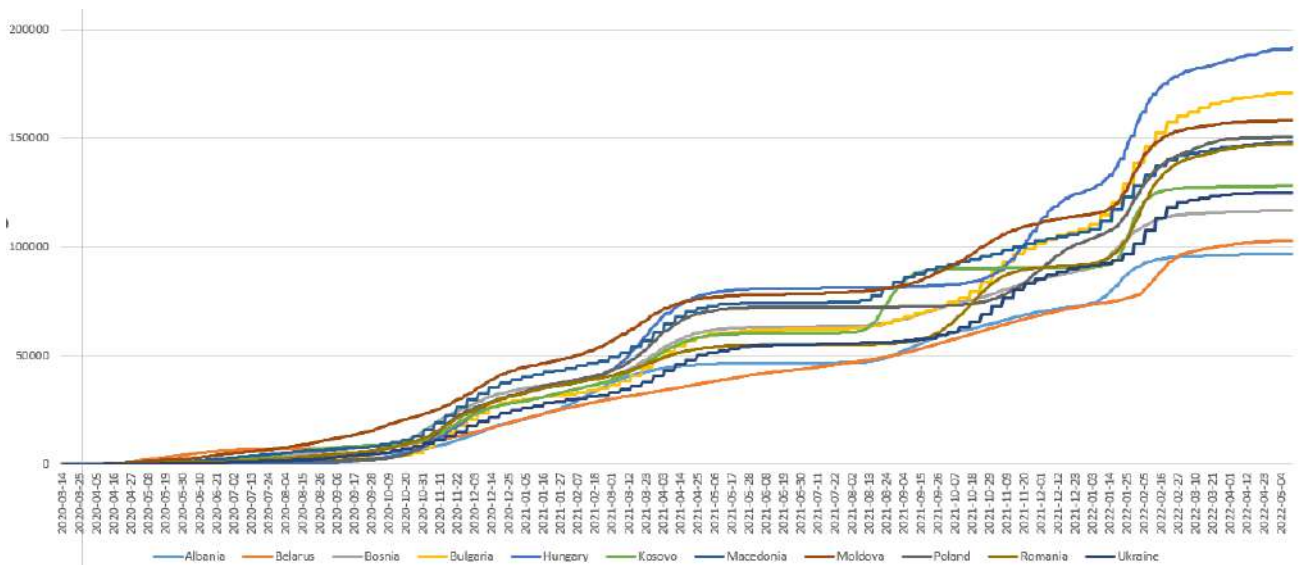


Рис. 19 Порівняльний графік загальної кількості випадків захворювання на мільйон населення для балканських та східноєвропейських країн з кластера, який був виявлений для кумулятивних даних

Додаток 2

У цьому додатку наведені таблиці 3 та 4.

Таблиця 3. Розмір найбільшої групи країн зі спільними кордонами в кластері для різних результатів кластеризації (початкові центроїди на основі результатів метода найближчого сусіда)

	Cluster 1	Cluster 2	Cluster 3	Cluster 4
Initial dataset, Euclidean metric	2	5	13	5
New dataset, Euclidean metric	2	2	8	7

New dataset, Mahalanobis distance	5	2	4	7
---	---	---	---	---

Таблиця 4. Розмір найбільшої групи країн зі спільними кордонами в кластері для різних результатів кластеризації (початкові центроїди на основі груп географічно близьких країн)

	Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5	Cluster 6	Cluster 7
Initial dataset, Euclidean metric	2	4	2	4	6	4	6
New dataset, Euclidean metric	1	2	1	2	2	5	7
New dataset, Mahalanobis distance	1	2	2	3	2	3	1

Додаток 3

У цьому додатку наведені програми в Python, за допомогою яких були реалізовані методи кластеризації.

1) Програма для метода найближчого сусіда та метода k-середніх з Евклідовою метрикою:

```
import pandas as pd
```

```
import numpy as np
```

```
import copy
```

```
import xlwings as xw
```

```
def clusterize_with_set_centers(df, centers, k):
```

```
    clusters = [[] for i in range(k)]
```

```
    countries = df.columns
```

```
    for country in countries:
```

```
        min_dist = np.linalg.norm(df[country] - centers[0])
```

```
        closest_center_num = 0
```

```
        for i in range(len(centers)):
```

```
            if np.linalg.norm(df[country] - centers[i]) < min_dist:
```

```
                min_dist = np.linalg.norm(df[country] - centers[i])
```

```
                closest_center_num = i
```

```
        clusters[closest_center_num].append(country)
```

```
    return clusters
```

```
def cluster_average(cluster):
```

```
    n = len(cluster)
```

```
    cluster_sum = sum([df[country] for country in cluster])
```

```
    return cluster_sum/n
```

```

def radius_and_distances(number, center, cluster):

    if len(cluster)==2:

        dist = [np.linalg.norm(df[country] - center) for country in cluster]

        print('Radius of the cluster', number, ': ', max(dist))

        print('Distances:')

        for i in range(2):

            print(cluster[i], ': ', dist[i])

        return dist

    else:

        dist = [np.linalg.norm(df[country] - center) for country in cluster]

        print('Radius of the cluster', number, ': ', max(dist))

        distance_dict = {key:value for (key,value) in zip(dist, cluster)}

        dist_sorted = copy.deepcopy(dist)

        dist_sorted.sort()

        print('Distances:')

        for d in dist_sorted:

            print(distance_dict[d], ': ', d)

        return dist

def neighbour_method(data, n):

    col = list(data.columns)

    last_col = col[-1]

    last_row = len(data)-1

    if n==0:

```

```

return

else:

    col = list(data.columns)

    dist = np.array([[np.linalg.norm(data.iloc[0:-1,i] - data.iloc[0:-1,j]) for j in range(0, len(col))] for i in
range(len(col))])

    mins = [min(dist[i, (i+1):len(col)]) for i in range(len(col)-1)]

    minimum = min(mins)

    min_indices = list(zip(*np.where(dist == minimum)))

    i_min = min_indices[0][0]

    j_min = min_indices[0][1]

    print('merge: ', data[i_min][last_row], 'and', data[j_min][last_row])

    print('distance: ', minimum)

    temp1 = data[last_col-1]

    temp2 = data[last_col]

    min1 = data[i_min]

    min2 = data[j_min]

    if (i_min == (last_col-1)) and (j_min != last_col):

        data[last_col] = min2

        data[j_min] = temp2

    if (i_min != (last_col-1)) and (j_min == last_col):

        data[last_col-1] = min1

```

```

    data[i_min] = temp1

if (j_min == last_col-1):

    data[last_col] = min1

    data[i_min] = temp2

if (i_min != (last_col-1)) and (j_min != last_col) and (j_min != (last_col-1)):

    data[last_col-1] = min1

    data[last_col] = min2

    data[i_min] = temp1

    data[j_min] = temp2

data.iloc[0:-1, last_col-1] = (data.iloc[0:-1, last_col-1]+data.iloc[0:-1, last_col])/2

data[last_col-1][last_row] = data[last_col-1][last_row] + '+' + data[last_col][last_row]

data1 = data.iloc[0:len(data), 0:last_col]

return neighbour_method(data1, n-1)

```

k = 4

```
wb = xw.Book('Final COVID data.xls')
```

```
data_excel = wb.sheets['new cases']
```

```
df = data_excel.range('C1:AP582').options(pd.DataFrame, header = True, index = False).value
```

```
print('NEIGHBOR CLUSTERS')
```

```
neighbor_clusters = [['UK', 'Ireland', 'Spain', 'Monaco', 'Portugal'], \
```

```
    ['Poland', 'Hungary', 'France', 'Netherland', 'Sweden'], \
```

```
    ['Austria', 'Italy', 'Bosnia', 'Bulgaria', 'Macedonia', 'Moldova', 'Romania'], \
```

```
    ['Germany', 'Ukraine', 'Russia', 'Belarus', 'Norway', 'Finland', 'Denmark']]
```

```
centers = [cluster_average(neighbor_clusters[i]) for i in range(k)]
```

```

clusters_old = clusterize_with_set_centers(df, centers, k)

while True:

    new_centers = [cluster_average(clusters_old[j]) for j in range(len(clusters_old))]

    clusters_new = clusterize_with_set_centers(df, new_centers, k)

    if clusters_new == clusters_old:

        for i in range(k):

            print('Cluster ', i+1, ': ', clusters_new[i])

            distances = radius_and_distances(i+1, new_centers[i], clusters_new[i])

            print(' ')

            pre_dist_df = [distances, clusters_new[i]]

            dist_df = pd.DataFrame(pre_dist_df, columns = [i for i in range(len(clusters_new[i]))])

            neighbour_method(dist_df, len(clusters_new[i])-1)

            print(' ')

        break

    else:

        clusters_old = copy.deepcopy(clusters_new)

```

```

k1 = 7

```

```

print('GEOGRAPHIC CLUSTERS')

```

```

geographic_clusters = [['Spain', 'Portugal'], \
                        ['France', 'Belgium', 'Netherland'], \
                        ['UK', 'Ireland'], \
                        ['Germany', 'Austria', 'Switzerland', 'Italy'], \

```

```

    ['Denmark', 'Sweden', 'Norway', 'Finland'], \
    ['Poland', 'Ukraine', 'Estonia', 'Latvia', 'Lithuania', 'Czech', 'Slovakia', 'Hungary'], \
    ['Slovenia', 'Croatia', 'Bosnia', 'Serbia', 'Macedonia', 'Albania', 'Greece', 'Romania', 'Moldova']]

centers_g = [cluster_average(geographic_clusters[i]) for i in range(k1)]

clusters_old_g = clusterize_with_set_centers(df, centers_g, k1)

while True:

    new_centers_g = [cluster_average(clusters_old_g[j]) for j in range(len(clusters_old_g))]

    clusters_new_g = clusterize_with_set_centers(df, new_centers_g, k1)

    if clusters_new_g == clusters_old_g:

        for i in range(k1):

            print('Cluster ', i+1, ': ', clusters_new_g[i])

            distances_g = radius_and_distances(i+1, new_centers_g[i], clusters_new_g[i])

            print(' ')

            pre_dist_df_g = [distances_g, clusters_new_g[i]]

            dist_df_g = pd.DataFrame(pre_dist_df_g, columns = [i for i in range(len(clusters_new_g[i]))])

            #print(dist_df)

            if len(clusters_new_g[i]) > 2:

                neighbour_method(dist_df_g, len(clusters_new_g[i])-1)

                print(' ')

            else:

                print("The cluster consists of two or less elements, therefore the nearest neighbor algorithm is
useless and non applicable")

                print(' ')

            break

    else:

```

```
clusters_old_g = copy.deepcopy(clusters_new_g)
```

2) Програма для метода k-середніх з метрикою Махаланобіса:

```
import pandas as pd

import numpy as np

import scipy

import copy

import xlwings as xw

from math import sqrt

wb = xw.Book('COVID data for Mahalanobis distance 1.xls')

data_excel = wb.sheets['new cases']

df = data_excel.range('D589:AQ593').options(pd.DataFrame, header = True, index = False).value

#print(df)

#print(df.to_numpy())

covariance_matrix = np.cov(df.to_numpy())

eps = 0

print('eps =', eps)

covariance_matrix = covariance_matrix + eps*np.identity(4)

covariance_matrix_inverse = np.linalg.inv(covariance_matrix)

#print(np.linalg.eigvals(covariance_matrix_inverse))

def mahalanobis(x, y, inv_cov):

    left_product = np.dot((x-y), inv_cov)

    #print('temp = ', np.dot(left_product, (x-y).T))
```

```

dist = sqrt(np.dot(left_product, (x-y).T))

return dist

def clusterize_with_set_centers(df, centers, k):

    clusters = [[] for i in range(k)]

    countries = df.columns

    for country in countries:

        min_dist = mahalanobis(df[country].to_numpy(), centers[0].to_numpy(), covariance_matrix_inverse)

        closest_center_num = 0

        for i in range(len(centers)):

            if mahalanobis(df[country].to_numpy(), centers[i].to_numpy(), covariance_matrix_inverse) <
min_dist:

                min_dist = mahalanobis(df[country].to_numpy(), centers[i].to_numpy(),
covariance_matrix_inverse)

                closest_center_num = i

            clusters[closest_center_num].append(country)

    return clusters

def cluster_average(cluster):

    n = len(cluster)

    cluster_sum = sum([df[country] for country in cluster])

    return cluster_sum/n

def radius_and_distances(number, center, cluster):

    if len(cluster)==2:

```

```
    dist = [mahalanobis(df[country].to_numpy(), center.to_numpy(), covariance_matrix_inverse) for
country in cluster]
```

```
    print('Radius of the cluster', number, ': ', max(dist))
```

```
    print('Distances:')
```

```
    for i in range(2):
```

```
        print(cluster[i], ': ', dist[i])
```

```
    return dist
```

```
else:
```

```
    dist = [mahalanobis(df[country].to_numpy(), center.to_numpy(), covariance_matrix_inverse) for
country in cluster]
```

```
    print('Radius of the cluster', number, ': ', max(dist))
```

```
    distance_dict = {key:value for (key,value) in zip(dist, cluster)}
```

```
    dist_sorted = copy.deepcopy(dist)
```

```
    dist_sorted.sort()
```

```
    print('Distances:')
```

```
    for d in dist_sorted:
```

```
        print(distance_dict[d], ': ', d)
```

```
    return dist
```

```
def neighbour_method(data, n):
```

```
    col = list(data.columns)
```

```
    last_col = col[-1]
```

```
    last_row = len(data)-1
```

```
    if n==0:
```

```
        return
```

```
    else:
```

```

col = list(data.columns)

#dist = np.array([[np.linalg.norm(data.iloc[0:-1,i] - data.iloc[0:-1,j]) for j in range(0, len(col))] for i in
range(len(col))])

dist = np.array([[abs(float(data.iloc[0:-1,i] - data.iloc[0:-1,j])) for j in range(0, len(col))] for i in
range(len(col))])

mins = [min(dist[i, (i+1):len(col)]) for i in range(len(col)-1)]

minimum = min(mins)

min_indices = list(zip(*np.where(dist == minimum)))

i_min = min_indices[0][0]
j_min = min_indices[0][1]

print('merge: ', data[i_min][last_row], 'and', data[j_min][last_row])

print('distance: ', minimum)

temp1 = data[last_col-1]
temp2 = data[last_col]

min1 = data[i_min]
min2 = data[j_min]

if (i_min == (last_col-1)) and (j_min != last_col):

    data[last_col] = min2

    data[j_min] = temp2

if (i_min != (last_col-1)) and (j_min == last_col):

    data[last_col-1] = min1

    data[i_min] = temp1

```

```

if (j_min == last_col-1):
    data[last_col] = min1
    data[i_min] = temp2
if (i_min != (last_col-1)) and (j_min != last_col) and (j_min != (last_col-1)):
    data[last_col-1] = min1
    data[last_col] = min2
    data[i_min] = temp1
    data[j_min] = temp2

data.iloc[0:-1, last_col-1] = (data.iloc[0:-1, last_col-1]+data.iloc[0:-1, last_col])/2
data[last_col-1][last_row] = data[last_col-1][last_row] + '+' + data[last_col][last_row]
data1 = data.iloc[0:len(data), 0:last_col]
return neighbour_method(data1, n-1)

'''
dist_test = mahalanobis(df['Albania'].to_numpy(), df['Bosnia'].to_numpy(), covariance_matrix_inverse)
print('dist_custom = ', dist_test)

print('dist_scipy = ', scipy.spatial.distance.mahalanobis(df['Albania'].to_numpy(), df['Bosnia'].to_numpy(),
covariance_matrix_inverse))
'''

k = 4

print('          MAHALANOBIS')

print('NEIGHBOR CLUSTERS')

```

```

neighbor_clusters = [['UK', 'Ireland', 'Spain', 'Monaco', 'Portugal'], \
                    ['Poland', 'Hungary', 'France', 'Netherland', 'Sweden'], \
                    ['Austria', 'Italy', 'Bosnia', 'Bulgaria', 'Macedonia', 'Moldova', 'Romania'], \
                    ['Germany', 'Ukraine', 'Russia', 'Belarus', 'Norway', 'Finland', 'Denmark']]

centers = [cluster_average(neighbor_clusters[i]) for i in range(k)]

clusters_old = clusterize_with_set_centers(df, centers, k)

while True:

    new_centers = [cluster_average(clusters_old[j]) for j in range(len(clusters_old))]

    clusters_new = clusterize_with_set_centers(df, new_centers, k)

    if clusters_new == clusters_old:

        for i in range(k):

            print('Cluster ', i+1, ': ', clusters_new[i])

            distances = radius_and_distances(i+1, new_centers[i], clusters_new[i])

            print(' ')

            pre_dist_df = [distances, clusters_new[i]]

            dist_df = pd.DataFrame(pre_dist_df, columns = [i for i in range(len(clusters_new[i]))])

            #print(dist_df)

            neighbour_method(dist_df, len(clusters_new[i])-1)

            print(' ')

        break

    else:

        clusters_old = copy.deepcopy(clusters_new)

```

k1 = 7

```

print('GEOGRAPHIC CLUSTERS')

geographic_clusters = [['Spain', 'Portugal'], \
                        ['France', 'Belgium', 'Netherland'], \
                        ['UK', 'Ireland'], \
                        ['Germany', 'Austria', 'Switzerland', 'Italy'], \
                        ['Denmark', 'Sweden', 'Norway', 'Finland'], \
                        ['Poland', 'Ukraine', 'Estonia', 'Latvia', 'Lithuania', 'Czech', 'Slovakia', 'Hungary'], \
                        ['Slovenia', 'Croatia', 'Bosnia', 'Serbia', 'Macedonia', 'Albania', 'Greece', 'Romania', 'Moldova']]

centers_g = [cluster_average(geographic_clusters[i]) for i in range(k1)]

clusters_old_g = clusterize_with_set_centers(df, centers_g, k1)

while True:

    new_centers_g = [cluster_average(clusters_old_g[j]) for j in range(len(clusters_old_g))]

    clusters_new_g = clusterize_with_set_centers(df, new_centers_g, k1)

    if clusters_new_g == clusters_old_g:

        for i in range(k1):

            print('Cluster ', i+1, ': ', clusters_new_g[i])

            distances_g = radius_and_distances(i+1, new_centers_g[i], clusters_new_g[i])

            print(' ')

            pre_dist_df_g = [distances_g, clusters_new_g[i]]

            dist_df_g = pd.DataFrame(pre_dist_df_g, columns = [i for i in range(len(clusters_new_g[i]))])

            #print(dist_df)

            if len(clusters_new_g[i]) > 2:

                neighbour_method(dist_df_g, len(clusters_new_g[i])-1)

                print(' ')

            else:

```

```
print("The cluster consists of two or less elements, therefore the nearest neighbor algorithm is  
useless and non applicable')
```

```
print('')
```

```
break
```

```
else:
```

```
clusters_old_g = copy.deepcopy(clusters_new_g)
```