

**Машинне
навчання**

Навчання з
прецедентами

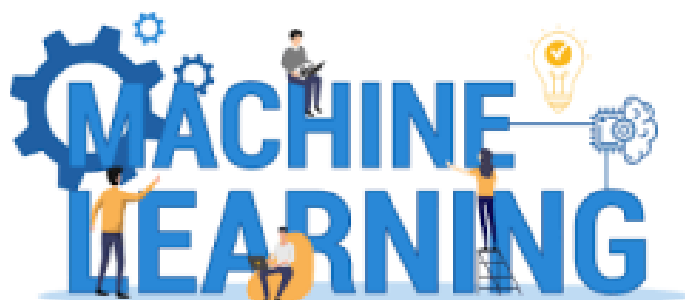
Методи
кластерного аналізу

Навчання з
підкріпленням

ЧАСТИНА 2

Методи кластерного аналізу

Антонюк С.В., Горбатенко М.Ю.,
Кириченко О.Л., Малик І.В.



Міністерство освіти і науки України
Чернівецький національний університет
імені Юрія Федьковича

Антонюк С.В., Горбатенко М.Ю., Кириченко О.Л., Малик І.В.

МАШИННЕ НАВЧАННЯ

Методи кластерного аналізу

Навчальний посібник

Чернівці
2023

УДК 004.8:519.7] (075.8)
М 382

Друкується за ухвалою Вченої ради Чернівецького національного університету імені Юрія Федьковича (протокол № 12 від 01.12.2021)

Рецензенти:

Сливка-Тилищак Ганна Іванівна, доктор фіз.-мат. наук, професор (Ужгородський національний університет);

Василик Ольга Іванівна, доктор фіз.-мат. наук, доцент (Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського»)

Машинне навчання. Методи кластерного аналізу: Навчальний посібник. / Антонюк С.В., Горбатенко М.Ю., Кириченко О.Л., Малик І.В. Чернівці : Чернів. нац. ун-т ім. Ю.Федьковича, 2023, 208 с.

Пропоноване видання містить теоретичний матеріал та модельні приклади з дисципліни «Методи та засоби кластерного аналізу». Подано основні методи і описано моделі машинного навчання без вчителя. Наведено велику кількість модельних прикладів, які дозволяють краще усвідомити теоретичний матеріал, а також сформовано завдання і методичні вказівки для лабораторних робіт по курсу. Навчальний посібник призначено для студентів технічних і фізико-математичних спеціальностей, які володіють базовими знаннями з математичного аналізу, теорії ймовірностей, математичної статистики, програмування. Для студентів навчально-наукового інституту фізико-технічних і комп'ютерних наук ЧНУ, а також студентів інших факультетів, які навчаються за спеціальностями «Комп'ютерні науки»

УДК 004.8 (519.7)

Зміст

1	Елементи чисельних методів	8
1.1	Метод Ньютона	8
1.1.1	Метод Ньютона для одновимірного випадку	8
1.1.2	Метод Ньютона для багатовимірного випадку	10
1.1.3	Модифікований метод Ньютона	13
1.2	Метод скінченних елементів	16
1.2.1	Е-М алгоритм	22
1.3	Модельовання випадкових величин	29
1.3.1	Модельовання одновимірних випадкових величин	29
1.3.2	Алгоритм Метрополіса - Гастінгса	31
1.3.3	Створення вибірок за Гіббсом	39
1.4	Вибір представників	43
1.5	Bootstrap метод	47
1.5.1	Непараметричний bootstrap	48
1.5.2	Байєсівський bootstrap	48
1.5.3	Згладжений bootstrap	49
1.5.4	Параметричний bootstrap	49
2	Алгоритми без вчителя. Кластеризація	51
2.1	Основи кластеризації	52

2.2	Основні типи змінних	53
2.3	Відстань та подібність	55
2.4	Теорема Клейнберга	59
2.5	Кластеризація розбиттям. Метод k -середніх	60
2.6	Ієрархічна кластеризація	67
2.6.1	Метод найближчого сусіда	68
2.6.2	Інші методи ієрархічної кластеризації	71
2.7	Складність алгоритмів	73
2.8	Нечітка кластеризація	73
2.8.1	Метод c -середніх	75
2.8.2	Метод Густафсона – Кесселя	78
2.8.3	Метод Газа – Гева	81
2.8.4	Моделі перемішування. Суміші.	82
2.9	Факторний аналіз. Метод головних компонент	91
2.9.1	Метод головних компонент	91
2.9.2	Вибір кількості головних компонент	97
2.9.3	Автокодування	99
2.9.4	Самоорганізаційні мапи Кохонена	101
2.9.5	Випадкове індексування	102
2.9.6	Приклад	103
2.10	Нормалізація факторів	107
2.11	Кластеризація нечислових даних	111
2.12	Проблема вибору оптимального числа кластерів	115
2.12.1	Метод зламаного тростини (метод ліктя)	115
2.12.2	Статистика розриву	117
2.12.3	Метод усередненого силуету	118
2.12.4	Інші методи вибору оптимальної кількості кластерів	120

2.13 Кластеризація на графах	123
2.13.1 Метод відсікання	127
2.13.2 Алгоритм Гірвана – Ньюмена	130
2.13.3 Марковський алгоритм кластеризації на гра- фах	133
2.13.4 PageRank алгоритм	138
2.14 Алгебраїчний підхід до вибору кількості кластерів	150
2.15 Кластеризація текстів	153
2.15.1 Попередня обробка тексту	154
2.15.2 Статистичні показники тексту	157
2.16 Кластеризація великих даних	160
2.16.1 Кластеризація на основі представників	161
2.16.2 BFR алгоритм	163
2.16.3 CURE алгоритм	169
2.16.4 BIRCH алгоритм	173
2.16.5 DBSCAN алгоритм	178
2.17 Кластеризація на основі процесу Діріхле	184
2.18 Статистика Хопкінса	196
3 Алгоритми з частковим навчанням	203
3.1 Основні гіпотези	205
3.2 Від класифікації до кластеризації	206
3.3 Клас трансдуктивних алгоритмів	208

Вступ

У даному підручнику буде розглянуто основні методи кластеризації або навчання без вчителя (unsupervised learning). Слід зауважити, що на відміну від навчання зі вчителем (supervised learning) в кластерному аналізі з'являються нові задачі, що пов'язані з визначенням оптимальної кількості кластерів k_{opt} , які не розглядаються в інших задачах машинного навчання. На відміну від задач класифікації, неможливо побудувати однозначної міри відповідності між знайденими кластерами та "реальними групами", оскільки даних груп (вчителя) не існує. У зв'язку з цим, в кластерному аналізі відсутні перевірки на основі точності та k-fold валідація.

Перший розділ присвячено деяким математичним методам, які будуть безпосередньо розглядатися при пошуку оптимальних параметрів в задачах кластеризації. Основна увага в цьому розділі буде приділена пошуку наближених розв'язків оптимізаційних задач та моделювання значень випадкових величин на основі заданого розподілу (щільності $f(x)$).

У другому розділі підручника розглянемо основні поняття кластерного аналізу, а саме, кластеру, матриці відстаней, матриці подібності, оптимального числа кластерів, сумішей, нечі-

ткої кластеризації. Також, увага буде приділено кластеризації нечислових даних (бінарних, номінальних та порядкових змінних). Поряд із цим, у цьому розділі будуть розглянуті основні методи кластеризації неструктурованих даних - кластеризація текстів та графів. Зауважимо, що методи які використовуються для кластеризації неструктурованих даних, відрізняються від класичних методів кластеризації числових даних оскільки неможливо однозначно визначити відстань чи подібність між об'єктами.

Четвертий розділ присвячений розгляду кластеризації у випадку часткової наявності вчителя, тобто за умови наявності y_i для деякого відсотку даних. Задача кластеризації в цьому випадку тісно переплетена із задачами класифікації, що і буде обговорено у пріоритетах розв'язання задач із частковим навчанням.

Термінологічний показчик

- BIRCH-алгоритм, 174;
- BFR-алгоритм, 164;
- Bootstrap-метод, 47;
- CURE-алгоритм, 169;
- DBSCAN-алгоритм, 179;
- EM-алгоритм, 22;
- PageRank-алгоритм, 139;
- Автокодування, 99;
- Алгоритм Гірвана – Ньюмена, 130;
- Бінарна змінна, 54;
- Відстань (метрика), 55;
- Випадкове індексування, 102;
- Інтервальна змінна, 54;
- Кластер, 52;

- Лема Джонсона - Лінденштрауса, 102;
- Коефіцієнт подібності, 57;
- Медоїд, 66;
- Метод Ньютона, 8;
- Метод усередненого силуету, 118;
- Метод k -середніх, 60;
- Номінальна змінна, 54;
- Об'єкт кластеризації, 52;
- Орієнтований граф, 123;
- Порядкова змінна, 54;
- Правило Кайзера, 98;
- Правило зламаної тростини, 99, 115;
- Процес Діріхле, 190;
- Самоорганізаційні мапи Кохонена, 101;
- Статистика Хопкінса, 197;
- Теорема Клейнберга, 59;
- Фактор об'єкта, 52;

Список літератури

- [1] Charu C. Aggarwal and Chandan K. Reddy, editors. *Data Clustering: Algorithms and Applications*. CRC Press, 2014.
- [2] B. Everitt, S. Landau, M. Leese, D. Stahl, and an O'Reilly Media Company Safari. *Cluster Analysis, 5th Edition*. John Wiley & Sons, 2011.
- [3] Leonard Kaufman and Peter J. Rousseeuw. *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley, 1990.
- [4] Jon M. Kleinberg. An impossibility theorem for clustering. pages 463–470, 2002.
- [5] O. Kyrychenko. Information technologies for statistical cluster analysis of information in complex networks. *Computer Systems and Information Technologies*, (4):47–51, 2022.
- [6] Thomas Reutterer and Daniel Dan. Cluster Analysis in Marketing Research. In Christian Homburg, Martin Klarmann, and Arnd Vomberg, editors, *Handbook of Market Research*, Springer Books, page 345. Springer, June 2022.

- [7] Peter Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.*, 20(1):53–65, 1987.
- [8] Остапов С.Е. Кириченко О.Л., Малик І.В. Стохастичні моделі в задачах штучного інтелекту. *Вісник Київського національного університету імені Тараса Шевченка. Серія фізико-математичні науки*, (2):53–57, 2021.
- [9] Остапов С.Е. Кириченко О.Л., Малик І.В. Аналіз кластерної структури Інтернет-мереж на основі випадкових матриць. *Проблеми керування та інформатики*, (1):37–46, 2022.

© Чернівецький національний
університет, 2023

© Антонюк С.В., Горбатенко М.Ю.,
Кириченко О.Л., Малик І.В., 2023

Антонюк Світлана Володимирівна;

Горбатенко Микола Юрійович;

Кириченко Оксана Леонідівна;

Малик Ігор Володимирович

Навчальне видання

МАШИННЕ НАВЧАННЯ

Методи кластерного аналізу

Навчальний посібник

Відповідальний за випуск Я.М. Дрінь

Літературний редактор О.В. Колодій

Підписано до друку 28.06.2023. Формат 60 x 84/16.
Папір офсетний. Друк різнографічний. Ум.друк арк. 4,8.
Обл.вид арк.4,0. Тираж 50. Зам. Н-222.
Друкарня Чернівецького національного університету
58012, Чернівці, вул. Кошобинського, 2

Свідоцтво суб'єкта видавничої справи ДК №891 від 08.04.2002 р.