

АНАЛІЗ ОСОБЛИВОСТЕЙ РЕАЛІЗАЦІЇ ЕМ-АЛГОРИТМУ ПРИ КЛАСТЕРИЗАЦІЇ СИСТЕМ СИГНАЛЬНИХ КОНСТРУКЦІЙ

Вступ

Одними з класів задач обробки сигналів та даних є задачі кластерного аналізу [1] та структурного прогнозування [2].

Задачі кластеризації полягають у алгоритмічній класифікації моделей, первинними даними для яких є результати спостережень, елементи даних та системи ознак, на групи (кластери).

Цей клас задач, як правило, відносять до складних комбінаторних задач, які характеризуються наявністю припущень та контекстністю розв'язуваної задачі, а також до задач статистичного розпізнавання образів, задач машинного навчання [1].

Задачі структурного прогнозування полягають у виявленні структурованості даних, вивченні функціональних залежностей між складовими елементами оброблюваних даних.

Цей клас задач відносять до задач машинного навчання [2].

Задачі кластерного аналізу та структурного прогнозування доповнюють одна одну при розв'язанні ряду практичних задач у галузі обробки сигналів та даних, пов'язаних з виявленням структур у досліджуваних системах даних та сигнально-кодових конструкцій, взаємозв'язків між ними, оцінюванні їх параметрів.

Одним з відомих та розповсюджуваних методів оцінювання параметрів моделей та кластеризації, який використовується при кластерному аналізі та структурному прогнозуванні, є *expectation-maximization* (EM) алгоритм [3–5].

Цей алгоритм та його модифікації використовуються для оцінювання параметрів компонент сумішей у вигляді гаусівських змішаних моделей, якими представляються досліджувані дані. EM-алгоритм також дозволяє здійснювати кластеризацію об'єктів, якими є елементи досліджуваних даних, на основі отримуваних апостеріор-

них оцінок імовірностей приналежності об'єктів кластерам (компонентам гаусівської змішаної моделі) [5–7].

Постановка проблеми

До особливостей реалізації EM-алгоритму при розв'язанні різних задач оцінювання параметрів, кластерного аналізу та структурного прогнозування, які суттєво можуть вплинути як на результат розв'язання задачі, так і на обчислювальну та алгоритмічну складність процесу розв'язання задачі, можна віднести такі:

1) чутливість EM-алгоритму до початкового наближення оцінюваних параметрів моделі, що може проявлятися у розв'язанні задачі у локальних екстремумах логарифмічної функції правдоподібності (гіпотези щодо параметрів моделі), максимізація якої є метою ітераційної процедури EM-алгоритму;

2) апріорна невизначеність кількості компонент суміші, наприклад гаусівської змішаної моделі, параметри якої оцінюються при реалізації EM-алгоритму, якщо у задачі кластеризації, яка розв'язується з його використанням, кількість кластерів є апріорі невідомою;

3) необхідність введення та використання додаткових критеріїв додавання або видалення компонент (кластерів) у модифікаціях EM-алгоритму з додаванням або видаленням компонент, від яких залежать ймовірності виникнення помилок 1-го та 2-го роду при кластеризації;

4) необхідність для багатьох задач враховувати контекст та фізичний зміст розв'язуваної задачі при інтерпретації отриманих розв'язків;

5) можливе виникнення при реалізації алгоритму невизначеностей (або сингулярностей), які підлягають розкриттю та контекстному поясненню відповідно до змісту розв'язуваної задачі з метою корекції та продовження реалізації ітераційної процедури EM-алгоритму.

Стаття присвячена аналізу особливостей реалізації EM-алгоритму, які можуть виникати при кластеризації систем сигнальних конструкцій та є пов'язаними з утворенням невизначеностей (на практиці ідентифікується програмно діленням на "0") у логарифмічній функції правдоподібності через виокремлення EM-алгоритмом порожніх кластерів, що унеможливорює алгоритмічне продовження реалізації подальших ітерацій EM-алгоритму до його збіжності.

Аналіз останніх досліджень і публікацій

До останніх досліджень та публікацій, в яких для розв'язання теоретичних та прикладних задач використовується EM-алгоритм та його модифікації, можна віднести дослідження у галузі розробки методів адаптивного прогнозування надійності технічних систем [8], метод багатосенсорного відстеження космічних об'єктів, у якому EM-алгоритм використовується для апроксимації статистичних характеристик у результатах спостережень [9], дослідження у галузі систем масового обслуговування, які пов'язані з розробкою методів оцінювання параметрів у системах типу $M/M/1/K$ з моделлю неоднорідного пуассонівського потоку, у яких EM-алгоритм використовується для отримання оцінок за методом максимальної правдоподібності [10] тощо. До типових ознак цих та інших підходів та методів, у яких використовується EM-алгоритм, можна віднести статистичний характер розв'язуваної задачі та отримання оцінок параметрів використовуваних для їх розв'язання статистичних моделей на основі принципу максимальної правдоподібності.

Постановка завдання

Метою статті є аналіз особливостей реалізації EM-алгоритму, які виникають через виокремлення EM-алгоритмом порожніх кластерів при розв'язанні задачі кластеризації систем сигнальних конструкцій та пов'язані з утворенням невизначеностей (математичних сингулярностей) у логарифмічній функції правдоподібності, максимізація якої проводиться під час ітераційної процедури EM-алгоритму.

Виклад основного матеріалу дослідження

Пояснимо проблемну ситуацію, яка зазначена у постановці завдання, на прикладі.

Нехай у системі обробки сигналів спостерігаються гармонічні сигнали $x_i(t)$, $i = \overline{1,10}$, які відрізняються один від одного фазами та можуть бути представлені моделлю (1).

$$x(t) = \sin\left(\frac{2\pi}{T}t + \Delta\varphi\right), \quad (1)$$

де T – період коливань;

$$\Delta\varphi_1 = -\frac{11}{180}\pi (-11^0), i = 1;$$

$$\Delta\varphi_2 = -\frac{3}{180}\pi (-3^0), i = 2;$$

$$\Delta\varphi_3 = +\frac{1}{180}\pi (+1^0), i = 3;$$

$$\Delta\varphi_4 = +\frac{6}{180}\pi (+6^0), i = 4;$$

$$\Delta\varphi_5 = +\frac{13}{180}\pi (+13^0), i = 5;$$

$$\begin{aligned} \Delta\varphi_6 &= +\frac{42}{180}\pi = -\frac{33}{180}\pi + \frac{75}{180}\pi = \\ &= 3\Delta\varphi_1 + \frac{75}{180}\pi (+42^0 = -33^0 + 75^0), i = 6; \end{aligned}$$

$$\begin{aligned} \Delta\varphi_7 &= +\frac{66}{180}\pi = -\frac{9}{180}\pi + \frac{75}{180}\pi = \\ &= 3\Delta\varphi_2 + \frac{75}{180}\pi (+66^0 = -9^0 + 75^0), i = 7; \end{aligned}$$

$$\begin{aligned} \Delta\varphi_8 &= +\frac{78}{180}\pi = +\frac{3}{180}\pi + \frac{75}{180}\pi = \\ &= 3\Delta\varphi_3 + \frac{75}{180}\pi (+78^0 = +3^0 + 75^0), i = 8; \end{aligned}$$

$$\begin{aligned} \Delta\varphi_9 &= +\frac{93}{180}\pi = +\frac{18}{180}\pi + \frac{75}{180}\pi = \\ &= 3\Delta\varphi_4 + \frac{75}{180}\pi (+93^0 = +18^0 + 75^0), i = 9; \end{aligned}$$

$$\begin{aligned} \Delta\varphi_{10} &= +\frac{114}{180}\pi = +\frac{39}{180}\pi + \frac{75}{180}\pi = \\ &= 3\Delta\varphi_5 + \frac{75}{180}\pi (+114^0 = +39^0 + 75^0), i = 10. \end{aligned}$$

Параметри сигналів у прикладі було обрано таким чином, щоб $x_i(t)$, $i = \overline{1,10}$, могли утворювати кластери $\{x_1(t), \dots, x_5(t)\}$ та $\{x_6(t), \dots, x_{10}(t)\}$ за ознакою фази, що також проілюстровано на рис. 1.

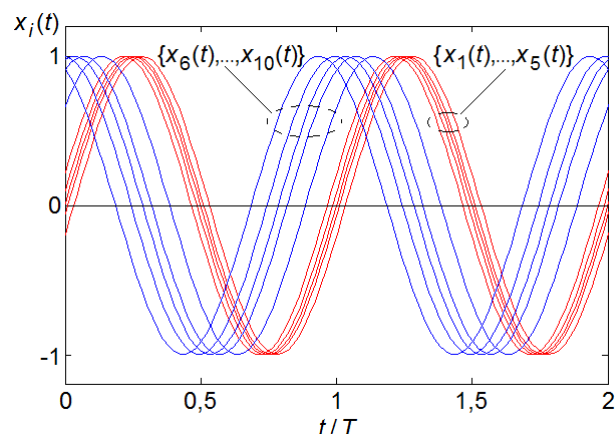


Рис. 1. Сигнали з прикладу, який розглядається

Під час аналізу сигналів при їх обробці в умовах апріорної невизначеності щодо структури та параметрів $x_i(t)$, $i = \overline{1,10}$, може бути поставлена задача кластеризації сигналів за такою ознакою кластеризації як взаємна кореляція між сигналами на спостережуваному інтервалі часу. Тоді вхідними даними для кластерного аналізу будуть значення нормованої взаємної кореляційної функції $r_{i,j}$, $i = \overline{1,10}$, $j = \overline{1,10}$, яка оцінюється, наприклад, на спостережуваному інтервалі часу $t \in (0, 2T)$ з припущенням, що математичне сподівання процесів $M\{x_i(t)\} = 0$:

$$r_{i,j} = \frac{1}{2\sigma\{x_i(t)\}\sigma\{x_j(t)\}T} \int_0^{2T} x_i(t)x_j(t)dt,$$

де $\sigma\{x_i(t)\}$, $\sigma\{x_j(t)\}$ — середньоквадратичні відхилення процесів $x_i(t)$ та $x_j(t)$ відповідно.

Для прикладу сигнальних конструкцій, який розглядається з використанням моделі (1), значення $r_{i,j}$, що містять інформацію про кореляції між усіма парами різних сигнальних конструкцій, можуть бути представлені таким чином:

$$\begin{aligned} \tilde{\mathbf{r}} = (\tilde{r}_n) = & (r_{1,2}; r_{1,3}; r_{1,4}; r_{1,5}; r_{1,6}; r_{1,7}; r_{1,8}; r_{1,9}; r_{1,10}; \\ & r_{2,3}; r_{2,4}; r_{2,5}; r_{2,6}; r_{2,7}; r_{2,8}; r_{2,9}; r_{2,10}; r_{3,4}; \\ & r_{3,5}; r_{3,6}; r_{3,7}; r_{3,8}; r_{3,9}; r_{3,10}; r_{4,5}; r_{4,6}; r_{4,7}; \\ & r_{4,8}; r_{4,9}; r_{4,10}; r_{5,6}; r_{5,7}; r_{5,8}; r_{5,9}; r_{5,10}; r_{6,7}; \\ & r_{6,8}; r_{6,9}; r_{6,10}; r_{7,8}; r_{7,9}; r_{7,10}; r_{8,9}; r_{8,10}; r_{9,10}), \\ & n = \overline{1, N}, N = 45. \end{aligned}$$

При вказаних вище параметрах моделі (1) вектор-рядок $\tilde{\mathbf{r}}$ матиме такий вигляд:

$$\begin{aligned} \tilde{\mathbf{r}} = & (0,990; 0,978; 0,956; 0,914; 0,602; \\ & 0,225; 0,017; -0,242; -0,574; 0,998; \\ & 0,988; 0,961; 0,707; 0,358; 0,156; \\ & -0,105; -0,454; 0,996; 0,978; 0,755; \\ & 0,423; 0,225; -0,035; -0,391; 0,993; \\ & 0,809; 0,500; 0,309; 0,052; -0,309; \\ & 0,875; 0,602; 0,423; 0,174; -0,191; \\ & 0,914; 0,809; 0,629; 0,309; 0,978; \\ & 0,891; 0,669; 0,966; 0,809; 0,934). \end{aligned} \quad (2)$$

Використання для кластеризації ЕМ-алгоритму полягає у статистичному аналізі $\tilde{\mathbf{r}}$ як даних, які відповідають певній суміші розподілів. Однією з моделей таких сумішей є гаусівська змішана модель [5–7], яка складається з гаусівсь-

ких компонент. Для випадку $\tilde{\mathbf{r}}$ ця модель формалізована виразом (3).

$$p(\tilde{r}) = \sum_{k=1}^K \frac{\theta_k}{\sqrt{2\pi\sigma_k^2}} \exp\left[-\frac{(\tilde{r} - m_k)^2}{2\sigma_k^2}\right], \quad (3)$$

де $p(\tilde{r})$ — щільність розподілу ймовірностей неперервної випадкової величини \tilde{r} ; K — кількість компонент суміші; θ_k — ваговий коефіцієнт k -ї гаусовської компоненти у складі суміші $p(\tilde{r})$; $\sum_{k=1}^K \theta_k = 1$; m_k — математичне сподівання k -ї гаусівської компоненти; σ_k — середньоквадратичне відхилення k -ї гаусівської компоненти; $k = \overline{1, K}$.

ЕМ-алгоритм дозволяє для спостережуваних даних $\tilde{\mathbf{r}}$ оцінити параметри гаусовської змішаної моделі $\omega = \{\theta, \mathbf{m}, \sigma\}$, а також апостеріорні ймовірності $\gamma_{n,k}$ приналежності елемента \tilde{r}_n k -ї компоненті суміші, які використовуються для прийняття рішення про приналежність елемента тому чи іншому кластеру.

Ітераційна процедура ЕМ-алгоритму [4, с.252] складається з кроку Е (expectation), на якому оцінюються значення $\gamma_{n,k}$ при поточному наближенні параметрів моделі, та кроку М (maximization), на якому уточнюються параметри моделі, максимізуючи у поточній ітерації логарифмічну функцію правдоподібності (гіпотези щодо параметрів моделі), яка формалізована для гаусівської змішаної моделі таким чином:

$$L(\omega | \tilde{\mathbf{r}}) = \sum_{n=1}^N \ln \sum_{k=1}^K \frac{\theta_k}{\sqrt{2\pi\sigma_k^2}} \exp\left[-\frac{(\tilde{r}_n - m_k)^2}{2\sigma_k^2}\right].$$

Формалізуємо процедуру ЕМ-алгоритму на його s -й ітерації ($s = 1, 2, 3, \dots$) для задачі кластеризації, яка розглядається у статті:

1) Е-крок: при наближенні параметрів моделі

$$\begin{aligned} \omega^{(s-1)} = & \{\theta^{(s-1)}, \mathbf{m}^{(s-1)}, \sigma^{(s-1)}\} = \\ = & (\theta_1^{(s-1)}, \dots, \theta_K^{(s-1)}, m_1^{(s-1)}, \dots, m_K^{(s-1)}, \sigma_1^{(s-1)}, \dots, \sigma_K^{(s-1)}) \end{aligned}$$

оцінюються значення

$$\gamma_{n,k}^{(s-1)} = \frac{\frac{\theta_k^{(s-1)}}{\sigma_k^{(s-1)}\sqrt{2\pi}} \exp\left[-\frac{(\tilde{r}_n - m_k^{(s-1)})^2}{2(\sigma_k^{(s-1)})^2}\right]}{\sum_{l=1}^K \frac{\theta_l^{(s-1)}}{\sigma_l^{(s-1)}\sqrt{2\pi}} \exp\left[-\frac{(\tilde{r}_n - m_l^{(s-1)})^2}{2(\sigma_l^{(s-1)})^2}\right]},$$

а також оцінюється кількість $N_k^{(s-1)}$ елементів $\tilde{\mathbf{r}}$, яка належить k -й компоненті суміші:

$$N_k^{(s-1)} = \sum_{n=1}^N \gamma_{n,k}^{(s-1)}. \quad (4)$$

У ітерації $s = 1$ використовується початкове наближення параметрів моделі:

$$\begin{aligned} \boldsymbol{\omega}^{(0)} &= \{\boldsymbol{\theta}^{(0)}, \mathbf{m}^{(0)}, \boldsymbol{\sigma}^{(0)}\} = \\ &= (\theta_1^{(0)}, \dots, \theta_K^{(0)}, m_1^{(0)}, \dots, m_K^{(0)}, \sigma_1^{(0)}, \dots, \sigma_K^{(0)}). \end{aligned}$$

2) М-крок: визначаються уточнені оцінки параметрів моделі

$$\begin{aligned} \boldsymbol{\omega}^{(s)} &= \{\boldsymbol{\theta}^{(s)}, \mathbf{m}^{(s)}, \boldsymbol{\sigma}^{(s)}\} = \\ &= (\theta_1^{(s)}, \dots, \theta_K^{(s)}, m_1^{(s)}, \dots, m_K^{(s)}, \sigma_1^{(s)}, \dots, \sigma_K^{(s)}), \end{aligned}$$

які у s -й ітерації максимізують $L(\boldsymbol{\omega} | \tilde{\mathbf{r}})$:

$$\begin{aligned} \theta_k^{(s)} &= \frac{N_k^{(s-1)}}{N}, \\ m_k^{(s)} &= \frac{1}{N_k^{(s-1)}} \sum_{n=1}^N \gamma_{n,k}^{(s-1)} \tilde{r}_n, \\ \sigma_k^{(s)} &= \sqrt{\frac{1}{N_k^{(s-1)}} \sum_{n=1}^N \gamma_{n,k}^{(s-1)} (\tilde{r}_n - m_k^{(s)})^2}. \end{aligned} \quad (5)$$

У випадку відсутності виникнення невизначеностей (математичних сингулярностей) у ітераційній процедурі алгоритму вона виконується до збіжності, практичним критерієм якої може бути умова $L(\boldsymbol{\omega}^{(s)} | \tilde{\mathbf{r}}) - L(\boldsymbol{\omega}^{(s-1)} | \tilde{\mathbf{r}}) < \varepsilon$, де $\varepsilon > 0$ — число, при якому точність оцінювання параметрів моделі вважається достатньою в контексті збіжності алгоритму.

Використовуватимемо у прикладі, який розглядається, такі параметри ($k = \overline{1, K}$):

- 1) $K = 11$;
- 2) $\theta_k^{(0)} = 1/K = 1/11$;

$$3) m_k^{(0)} = -1 + (2k - 1)/K = -1 + (2k - 1)/11;$$

$$4) \sigma_k^{(0)} = 1/3K = 1/33.$$

Вибір кількості компонент гаусівської змішаної моделі $K = 11$ було обрано таким чином, щоб вона мала достатньо велике значення у порівнянні з об'ємом вибірки ($N = 45$), що спричинить виявлення у ітераційній процедурі ЕМ-алгоритму принаймні одного порожнього кластеру та відповідної невизначеності (математичної сингулярності), яка робить неможливою подальшу алгоритмічну реалізацію ітераційної процедури ЕМ-алгоритму.

Вибір початкового наближення параметрів моделі $\boldsymbol{\omega}^{(0)}$ було обрано так, щоб вагові коефіцієнти $\theta_k^{(0)}$ були апіорі однаковими, у діапазоні можливих значень нормованої взаємної кореляційної функції розподіл математичних сподівань $m_k^{(0)}$ гаусовських компонент був рівномірним, а також використовувалося правило $\pm 3\sigma_k^{(0)}$ для їх середньоквадратичних відхилень.

Реалізація ЕМ-алгоритму для даних $\tilde{\mathbf{r}}$, які вказано у (2), а також при зазначених вище K та $\boldsymbol{\omega}^{(0)}$, дає такі результати процесу максимізації $L(\boldsymbol{\omega} | \tilde{\mathbf{r}})$ у ітераціях: $L(\boldsymbol{\omega}^{(0)} | \tilde{\mathbf{r}}) = -68,9$, $L(\boldsymbol{\omega}^{(1)} | \tilde{\mathbf{r}}) = -8,64$, $L(\boldsymbol{\omega}^{(2)} | \tilde{\mathbf{r}}) = 0/0$.

Виникнення невизначеності (математичної сингулярності) у $L(\boldsymbol{\omega}^{(2)} | \tilde{\mathbf{r}})$ пов'язане з тим, що утворюється порожній кластер $k = 1$, для якого виконується умова $\forall n \in \overline{1, N}, k = 1, \gamma_{n,k}^{(1)} \rightarrow 0$.

Проаналізуємо цю невизначеність для загального випадку утворення у s -й ітерації k -го порожнього кластеру, тобто $\lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} L(\boldsymbol{\omega}^{(s)} | \tilde{\mathbf{r}})$.

Представляючи $\boldsymbol{\omega}^{(s)}$ через $\boldsymbol{\gamma}^{(s-1)}$ з використанням (4) та (5), отримаємо:

$$\begin{aligned} \lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} L(\boldsymbol{\omega}^{(s)} | \tilde{\mathbf{r}}) &= \lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \sum_{n=1}^N \ln \sum_{k=1}^K \frac{\theta_k^{(s)}}{\sqrt{2\pi(\sigma_k^{(s)})^2}} \exp \left[-\frac{(\tilde{r}_n - m_k)^2}{2(\sigma_k^{(s)})^2} \right] \equiv \\ &\equiv \lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \sum_{n=1}^N \ln \sum_{k=1}^K \frac{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} \right)^{\frac{3}{2}}}{N \sqrt{2\pi \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2}} \exp \left[-\frac{(\tilde{r}_n - m_k^{(s)})^2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)}}{2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2} \right]. \end{aligned}$$

Проаналізуємо $A_{n,k} = \lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)}\right)^{\frac{3}{2}}}{N \sqrt{2\pi \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2}} \exp \left[-\frac{(\tilde{r}_n - m_k^{(s)})^2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)}}{2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2} \right]$. Для зручності ви-

користаємо логарифмування та властивість границі неперервної функції, тобто знайдемо

$$A_{n,k} = \exp(\ln A_{n,k}) = \exp \left[\lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \ln \left\{ \frac{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)}\right)^{\frac{3}{2}}}{N \sqrt{2\pi \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2}} \exp \left[-\frac{(\tilde{r}_n - m_k^{(s)})^2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)}}{2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2} \right] \right\} \right] =$$

$$= \exp \left[\lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \left(\ln \frac{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)}\right)^{\frac{3}{2}}}{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2\right)^{\frac{1}{2}}} - \ln N - \frac{1}{2} \ln(2\pi) - \frac{(\tilde{r}_n - m_k^{(s)})^2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)}}{2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2} \right) \right]$$

Для $\lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)}\right)^{\frac{3}{2}}}{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2\right)^{\frac{1}{2}}}$ використаємо правило Лопітала, враховуючи також те, що для

порожнього k -го кластеру $\tilde{r}_n - m_k^{(s)} \neq 0$:

$$\lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)}\right)^{\frac{3}{2}}}{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2\right)^{\frac{1}{2}}} = \lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{\frac{\partial}{\partial \gamma_{n,k}^{(s-1)}} \left\{ \left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)}\right)^{\frac{3}{2}} \right\}}{\frac{\partial}{\partial \gamma_{n,k}^{(s-1)}} \left\{ \left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2\right)^{\frac{1}{2}} \right\}} =$$

$$= \lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{\frac{3}{2} \left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)}\right)^{\frac{1}{2}}}{\frac{1}{2} (\tilde{r}_n - m_k^{(s)})^2 \left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2\right)^{\frac{1}{2}}} = \lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{3 \left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)}\right)^{\frac{1}{2}} \left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2\right)^{\frac{1}{2}}}{(\tilde{r}_n - m_k^{(s)})^2} = 0.$$

Для $\lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{(\tilde{r}_n - m_k^{(s)})^2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)}}{2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2}$ також використаємо правило Лопітала:

$$\lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{(\tilde{r}_n - m_k^{(s)})^2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)}}{2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2} = \lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, N}} \frac{\frac{\partial}{\partial \gamma_{n,k}^{(s-1)}} \left\{ (\tilde{r}_n - m_k^{(s)})^2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} \right\}}{\frac{\partial}{\partial \gamma_{n,k}^{(s-1)}} \left\{ 2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2 \right\}} = \frac{(\tilde{r}_n - m_k^{(s)})^2}{2 (\tilde{r}_n - m_k^{(s)})^2} = \frac{1}{2}.$$

З урахуванням цього

$$\lim_{\substack{\gamma_{n,k}^{(s-1)} \rightarrow 0, \\ n=1, \bar{N}}} \left(\ln \frac{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} \right)^{\frac{3}{2}}}{\left(\sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2 \right)^{\frac{1}{2}}} - \ln N - \frac{1}{2} \ln(2\pi) - \frac{(\tilde{r}_n - m_k^{(s)})^2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)}}{2 \sum_{p=1}^N \gamma_{p,k}^{(s-1)} (\tilde{r}_p - m_k^{(s)})^2} \right) = -\infty,$$

$A_{n,k} = \exp(\ln A_{n,k}) = \exp(-\infty) = 0$, $n = \overline{1, \bar{N}}$, для окремо взятого k -го порожнього кластеру.

Таким чином, у випадку утворення в ітераційній процедурі EM-алгоритму порожнього k -го кластеру та відповідних невизначеностей (математичних сингулярностей) для їх усунення необхідне виключення відповідної k -ї виродженої гаусівської компоненти на s -й ітерації EM-алгоритму зі складу $L(\omega^{(s)} | \tilde{\mathbf{r}})$, що не вплине на значення $L(\omega^{(s)} | \tilde{\mathbf{r}})$ через те, що для окремо взятого порожнього k -го кластеру $A_{n,k} = 0$, $n = \overline{1, \bar{N}}$.

Зауважимо, що у цій статті не наводиться повне розв'язання задачі кластеризації для використаного прикладу аналізу сигнальних конструкцій, а він служить виключно для пояснення можливого випадку виникнення математичної сингулярності в ітераційній процедурі EM-алгоритму через утворення порожніх кластерів з метою її аналізу. Також зазначимо, що в ітераційній процедурі EM-алгоритму можуть виникати невизначеності (математичні сингулярності) інших типів, які не пов'язані з утворенням порожніх кластерів та аналіз яких потребує окремого розгляду.

Висновки

Однією з особливостей реалізації EM-алгоритму при розв'язанні задач кластеризації систем сигнальних конструкцій є можливе виникнення стану математичної сингулярності у його ітераційній процедурі.

У статті показано приклад можливої задачі кластеризації сигнальних конструкцій, підхід до розв'язання якої полягає у аналізі взаємнокореляційних зв'язків між сигналами з використанням гаусівської змішаної моделі та оцінювання її параметрів і прихованих змінних (імовірностей приналежності елементів суміші певним компонентам цієї суміші, що є визначальним при прийнятті рішення про приналежність того чи іншого елемента певному кластеру) з використанням EM-алгоритму.

Розглянута у статті задача характеризується математичною сингулярністю, яка пов'язана з утворенням порожніх кластерів. У результаті аналізу математичної сингулярності такого типу (невизначеність "0/0") показано, що можливе

видалення утвореного порожнього кластеру таким чином, що структура та значення логарифмічної функції правдоподібності корегуються до таких, які були б у випадку апріорної відсутності зазначеного порожнього кластеру.

Проаналізована у статті особливість є більш властивою для модифікації EM-алгоритму з видаленням компонент, а також у випадках відносно невеликої кількості елементів, які підлягають кластеризації, або відносно великої кількості кластерів, значення якої є структурним параметром EM-алгоритму.

Залежно від задачі кластеризації, розглянута особливість EM-алгоритму може бути використана при кластерному аналізі складних конструкцій сигналів та даних, наприклад мультиплікативно комплементарних бінарних сигнально-кодових конструкцій [11], структур їх кореляційних функцій [12], систем корельованих сигнальних конструкцій при неортогональному множинному доступі (NOMA) у системах мобільних телекомунікацій 5G [13] тощо.

ЛІТЕРАТУРА

1. **Jain A. K.**, Murty M. N., Flynn P. J. Data clustering: a review. *ACM Computing Surveys*. 1999. Vol. 31. No. 3. P. 264-323. DOI: 10.1145/331499.331504.
2. **Bakir G. H.** et al. Predicting Structured Data (Neural Information Processing series) 1st Edition. The MIT Press, 2007. 360 p.
3. **Dempster A. P.**, Laird N. M., Rubin D. B. Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society. Series B*. 1977. Vol. 39. No. 1. P. 1-38.
4. **Gupta M. R.**, Chen Y. Theory and Use of the EM Algorithm. *Foundations and Trends® in Signal Processing*. 2011. Vol. 4. No. 3. P. 223-296. DOI: 10.1561/20000000034.
5. **Vlassis N.**, Likas A. A Greedy EM Algorithm for Gaussian Mixture Learning. *Neural Processing Letters*. 2002. Vol. 15. P. 77-87.
6. **Yu D.**, Deng L. Gaussian Mixture Models. *Automatic Speech Recognition. A Deep Learning Approach*. London: Springer-Verlag, 2015. P. 13-21. DOI: 10.1007/978-1-4471-5779-3_2

7. **Huang T.**, Peng H., Zhang K. Model Selection for Gaussian Mixture Models. *Statistica Sinica*. 2017. Vol. 27. P. 147-169.

8. **Wang Xi et al.** An Adaptive Prognostic Approach for Newly Developed System with Three-Source Variability. *IEEE Access*. 2019. Vol. 7. P. 53091-53102. DOI: 10.1109/ACCESS.2019.2911307

9. **Wei B.**, Nener B. D. Multi-Sensor Space Debris Tracking for Space Situational Awareness with Labeled Random Finite Sets. *IEEE Access*. 2019. Vol. 7. P. 36991-37003. DOI: 10.1109/ACCESS.2019.2904545

10. **Li C.**, Okamura H., Dohi T. Parameter Estimation of $M_r/M/1/K$ Queueing Systems with Utilization Data. *IEEE Access*. 2019. Vol. 7. P. 42664-42671. DOI: 10.1109/ACCESS.2019.2906796

11. **Голубничий А. Г.**, Коначович Г. Ф. Мультипликативно комплементарные бинарные сигнално-кодовые конструкции. *Известия высших учебных заведений. Радиоэлектроника*. 2018. Т. 61. № 10. С. 551-565. DOI: 10.20535/S0021347018100011.

12. **Голубничий О. Г.** Синтез аналітичних форм опису автокореляційної функції узагальнених бінарних послідовностей Баркера типу 1 на основі її декомпозиції з використанням лінійних складових. *Наукоємні технології*. 2019. Т. 41. № 1. С. 10–15. DOI: 10.18372/2310-5461.41.13523.

13. **Голубничий О. Г.** Синтез систем корельованих сигналів з використанням доповненої процедури Грама-Шмідта. *Наукоємні технології*. 2018. Т. 40. № 4. С. 405-408. DOI: 10.18372/2310-5461.40.13265.

Голубничий О. Г.

АНАЛІЗ ОСОБЛИВОСТЕЙ РЕАЛІЗАЦІЇ ЕМ-АЛГОРИТМУ ПРИ КЛАСТЕРИЗАЦІЇ СИСТЕМ СИГНАЛЬНИХ КОНСТРУКЦІЙ

EM-алгоритм (expectation-maximization algorithm) є відомим статистичним методом, який використовується в області обробки даних та сигналів для кластерного аналізу, оцінювання параметрів, а також в інших методах машинного навчання. EM-алгоритм характеризується деякими специфічними особливостями, наприклад чутливістю до початкових параметрів. Метою статті є аналіз особливостей EM-алгоритму, які виникають внаслідок виявлення порожніх кластерів при розв'язанні задачі кластеризації сигналних конструкцій. Ці особливості проявляються в утворенні невизначеностей (математичних сингулярностей) у логарифмічній функції правдоподібності, максимізація якої виконується під час ітеративної процедури EM-алгоритму. У статті наведено приклад можливої задачі кластеризації в області аналізу сигналів. У цьому прикладі використовується підхід, що базується на аналізі взаємних кореляцій між сигналами при представленні цих кореляцій гаусівською змішаною моделлю, а також на оцінюванні параметрів гаусівської змішаної моделі та прихованих змінних (ймовірностей приналежності елементів суміші до певних компонент цієї суміші, що є визначальним при прийнятті рішень про приналежність конкретного елемента до певного кластеру) за допомогою EM-алгоритму. В результаті аналізу розглянутого типу математичної сингулярності показано, що можна видалити порожній кластер таким чином, що структура і значення логарифмічної функції правдоподібності коригуються до тих, які були б у випадку апріорної відсутності зазначеного порожнього кластера. Особливість, яка проаналізована у статті, більш характерна для модифікації EM-алгоритму з видаленням компонент гаусівської змішаної моделі. Ця особливість може також виникати у випадках відносно невеликого числа елементів, які підлягають кластерному аналізу, або у випадку відносно великої кількості компонент (кластерів) гаусівської змішаної моделі, яка є структурним параметром EM-алгоритму.

Ключові слова: кластерний аналіз; EM-алгоритм; математична сингулярність; гаусівська змішана модель; машинне навчання; обробка даних.

Holubnychi A. G.

ANALYSIS OF PECULIARITIES OF THE EM-ALGORITHM IN THE IMPLEMENTATION OF CLUSTERING OF SIGNAL CONSTRUCTIONS

The expectation-maximization (EM) algorithm is a well-known statistical method, which is used in the field of data and signal processing for cluster analysis, parameter estimation and other machine learning techniques. The EM-algorithm is characterized by some specific peculiarities, e.g., a sensitivity to initial parameters. The article aims to analyze peculiarities of the EM-algorithm, which arise due to the detection of empty clusters in solving the problem of clustering of signal constructions. These peculiarities boil down to the formation of uncertainties (mathematical singularities) in the log-likelihood function, the maximization of which is performed during the iterative procedure of the EM-algorithm. An example of a possible clustering problem in the field of signal analysis is shown. The example uses an approach based on an analysis of correlations between signals using the Gaussian mixture model for these correlations and an estimation of the Gaussian mixture model parameters and hidden variables (the probabilities for belonging of elements of the mixture to certain components of this mixture, which is decisive for a criteria whether a particular element belongs to a particular cluster) using the EM-algorithm. As a result of the analysis of considered type of mathematical singularity, it is shown that it is possible to remove an empty cluster in such a way that the structure and

values of the log-likelihood function are adjusted to those that would be in the case of the a priori absence of this specified empty cluster. The peculiarity, which is analyzed in the article, is more typical for the modification of the EM-algorithm with deletion the components of the Gaussian mixture model. This peculiarity may also arise in cases of relatively small number of elements that need to be analyzed through clustering or relatively large number of components of the Gaussian mixture model (clusters), which is the structural parameter of the EM-algorithm.

Keywords: cluster analysis; expectation-maximization algorithm; mathematical singularity; Gaussian mixture model; machine learning; data processing.

Голубничий А. Г.

АНАЛИЗ ОСОБЕННОСТЕЙ РЕАЛИЗАЦИИ EM-АЛГОРИТМА ПРИ КЛАСТЕРИЗАЦИИ СИСТЕМ СИГНАЛЬНЫХ КОНСТРУКЦИЙ

EM-алгоритм (expectation-maximization algorithm) является известным статистическим методом, который используется в области обработки данных и сигналов в целях кластерного анализа, оценки параметров, а также в других методах машинного обучения. EM-алгоритм характеризуется некоторыми специфическими особенностями, например, чувствительностью к начальным параметрам. Целью статьи является анализ особенностей EM-алгоритма, которые возникают в результате обнаружения пустых кластеров при решении задачи кластеризации сигнальных конструкций. Эти особенности проявляются в образовании неопределенностей (математических сингулярностей) в логарифмической функции правдоподобия, максимизация которой выполняется при итерационной процедуре EM-алгоритма. В статье приведен пример возможной задачи кластеризации в области анализа сигналов. В этом примере используется подход, основанный на анализе взаимных корреляций между сигналами при представлении этих корреляций гауссовской смешанной моделью, а также на оценивании параметров гауссовской смешанной модели и скрытых переменных (вероятностей принадлежности элементов смеси определенным компонентам этой смеси, которые являются определяющими при принятии решений о принадлежности конкретного элемента к определенному кластеру) с помощью EM-алгоритма. В результате анализа рассмотренного типа математической сингулярности показано, что является возможным удаление пустого кластера таким образом, что структура и значение логарифмической функции правдоподобия корректируются так, если априори отсутствовал бы указанный пустой кластер. Особенность, которая проанализирована в статье, более характерна для модификации EM-алгоритма с удалением компонент гауссовской смешанной модели. Эта особенность может также возникать в случаях относительно небольшого числа элементов, подлежащих кластерному анализу, или в случае относительно большого количества компонент (кластеров) гауссовской смешанной модели, которое является структурным параметром EM-алгоритма.

Ключевые слова: кластерный анализ; EM-алгоритм; математическая сингулярность; гауссовская смешанная модель; машинное обучение; обработка данных.

Стаття надійшла до редакції 28.04.2019 р.

Прийнято до друку 31.05.2019 р.