

Стукалов А. С.

МГУ им. М. В. Ломоносова, ast@pochta.ru

ОЦЕНКА КАЧЕСТВА КЛАСТЕРИЗАЦИИ ВЫРОВНЕННЫХ СТРОК

Вопрос нахождения кластеров в наборе строк одинаковой длины является одним из актуальных вопросов биоинформатики. Вводится векторный критерий, позволяющий оценивать качество таких кластеризаций, рассматриваются способы его свертки. Данный критерий может быть использован для построения методов кластеризации.

1. Введение

Кластеризация данных [1], т. е. разбиение набора данных на подмножества, каждое из которых обладает некоторым специфическим для его элементов свойством, находит много применений в биоинформатике, экономике, распознавании образов и т. д. [2]

В частности, в биоинформатике кластеризация используется для выделения семейств гомологичных белков [3, 5, 6]. Обычно входные данные — это *множественное выравнивание (multiple sequence alignment, MSA)*, набор выровненных друг относительно друга белковых последовательностей. MSA представляет собой матрицу, каждая строка которой есть последовательность аминокислотных остатков некоторого белка. В эту строку вставлены специальные *символы-пробелы* таким образом, чтобы каждая колонка матрицы содержала преимущественно одинаковые или похожие остатки.

Одна из особенностей матрицы MSA состоит в том, что кроме *консервативных* колонок, состоящих из похожих аминокислотных остатков, в ней есть колонки, содержащие остатки существенно различной природы. Наличие таких колонок объясняется двумя причинами. Это, во-первых, функциональные и структурные различия белков, образующих матрицу. Во-вторых, *молчащие мутации*, которые не приводят к изменению структуры или функций белка.

Идеальный метод кластеризации белков должен уметь отличать колонки MSA, которые коррелированы со свойствами кластеров (*специфические колонки*) от неинформативных колонок, в которых преобладают

молчащие мутации. Специфические колонки должны быть ключевыми признаками, по которым ведется кластеризация: последовательности, содержащие на специфических позициях одинаковые остатки, должны попадать в один кластер. Кластеры, полученные таким способом, будут однозначно характеризоваться набором остатков на специфических позициях.

Можно задаться вопросом оценки качества кластеризации, введения некоторой меры, которая бы оценивала свойства, требуемые от хорошей кластеризации. Подобная мера могла бы быть использована в методе кластеризации, основанном на оптимизации некоторого функционала на множестве кластеризаций. Одним из примеров такого подхода является метод оптимизации функции контраста энтропии [5].

В данной статье предлагается двухкомпонентный векторный критерий оценки качества кластеризации. Его первый компонент оценивает *неоднородность строк*, составляющих кластер. Второй компонент оценивает *неоднозначность классифицирования*, возникающую из-за похожести строк разных кластеров друг на друга.

2. Постановка задачи

Пусть A_1, A_2, \dots, A_M — набор слов длины N в алфавите $\Lambda = \{\lambda_1, \lambda_2, \dots, \lambda_L\}$. Эти слова можно считать строками матрицы $A \in \Lambda^{M \times N}$. Пусть \tilde{A} — разбиение множества строк, т. е. $\tilde{A} = \{A_1, A_2, \dots, A_m\}$, где

$$A_i = \{A_{i_1}, A_{i_2}, \dots, A_{i_{n_i}}\} \neq \emptyset, \quad \bigsqcup_{A_i \in \tilde{A}} A_i = \{A_1, A_2, \dots, A_M\}. \quad (1)$$

Пусть \mathfrak{A} — множество всех таких разбиений \tilde{A} .

Предположим, существует мера *качества кластеризации*, или стоимости кластеризации, $q: \mathfrak{A} \rightarrow \mathbb{R}$. Задача состоит в нахождении оптимального разбиения \tilde{A}^* :

$$q(\tilde{A}^*) \leq q(\tilde{A}) \quad \forall \tilde{A} \in \mathfrak{A}. \quad (2)$$

В следующих пунктах будет рассмотрен один из возможных способов задания функции $q(\tilde{A})$.

3. Мера неоднородности состава

Пусть на множестве пар символов алфавита Λ определена функция $\rho: \Lambda \times \Lambda \rightarrow \mathbb{R}^+$, которая имеет смысл расстояния между двумя символами: чем больше различие между α и $\beta \in \Lambda$, тем больше $\rho(\alpha, \beta)$ (однако выполнение тождества $\rho(\alpha, \alpha) = 0$ не требуется). Простейший пример такой ρ — это *дискретная метрика* на Λ :

$$\rho(\alpha, \beta) = \begin{cases} 0, & \alpha = \beta; \\ 1, & \alpha \neq \beta. \end{cases} \quad (3)$$

Для задания расстояния в кластеризациях MSA могут быть использованы также *матрицы замен*, такие как BLOSUM или PAM [4], отражающие наблюдаемые частоты мутаций аминокислотных остатков.

Рассмотрим символы на j -й позиции в подмножестве строк $\mathcal{A} = \{A_{i_1}, A_{i_2}, \dots, A_{i_{n_{\mathcal{A}}}}\}$. Пусть $n_{\mathcal{A}}(\lambda, j)$ — количество появлений $\lambda \in \Lambda$ в j -й колонке \mathcal{A} .

Определение 1. Неоднородность состава множества строк \mathcal{A} в j -й колонке есть

$$d(\mathcal{A}, j) = \frac{1}{n_{\mathcal{A}}} \sum_{\alpha, \beta \in \Lambda} \rho(\alpha, \beta) n_{\mathcal{A}}(\alpha, j) n_{\mathcal{A}}(\beta, j). \quad (4)$$

Каждой паре (\mathcal{A}, j) мы можем сопоставить функцию вероятности

$$p_{\mathcal{A}, j}(\lambda) = \frac{1}{n_{\mathcal{A}}} n_{\mathcal{A}}(\lambda, j),$$

тогда $d(\mathcal{A}, j)$ можно переписать так:

$$d(\mathcal{A}, j) = n_{\mathcal{A}} \mathbb{E} \rho(\Lambda_1, \Lambda_2), \quad (5)$$

где Λ_1 и Λ_2 — две независимые случайные величины, принимающие значения из Λ , с функцией вероятности $p_{\mathcal{A}, j}$.

Рассмотрим случай, когда ρ — это *дискретная метрика* на Λ , тогда

$$\mathbb{E} \rho(\Lambda_1, \Lambda_2) = p(\Lambda_1 \neq \Lambda_2).$$

Минимум $d(\mathcal{A}, j) = 0$ достигается, когда в каждой строке из \mathcal{A} j -м символом является некоторый $\lambda \in \Lambda$. Напротив, если в j -й колонке \mathcal{A} каждый символ встречается единожды, неоднородность состава достигает максимального значения $d(\mathcal{A}, j) = n_{\mathcal{A}} - 1$.

Пусть $\tilde{\mathcal{A}} = \{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_{n_{\tilde{\mathcal{A}}}}\}$ — попарно непересекающиеся подмножества строк \mathcal{A} . Мы расширяем определение меры неоднородности строк и полагаем

$$d(\tilde{\mathcal{A}}, j) = \sum_{i=1}^{n_{\tilde{\mathcal{A}}}} d(\mathcal{A}_i, j). \quad (6)$$

Интересно установить, как разделение и объединение кластеров влияют на эту меру, т. е. сравнить $d(\mathcal{A}, j) + d(\mathcal{B}, j)$ и $d(\mathcal{A} \cup \mathcal{B}, j)$:

$$\begin{aligned} d(\mathcal{A}) + d(\mathcal{B}) - d(\mathcal{A} \cup \mathcal{B}) &= \\ &= \sum_{\alpha, \beta \in \Lambda} \rho(\alpha, \beta) \frac{n_{\mathcal{A}}(\alpha) n_{\mathcal{A}}(\beta)}{n_{\mathcal{A}}} + \sum_{\alpha, \beta \in \Lambda} \rho(\alpha, \beta) \frac{n_{\mathcal{B}}(\alpha) n_{\mathcal{B}}(\beta)}{n_{\mathcal{B}}} - \\ &- \sum_{\alpha, \beta \in \Lambda} \rho(\alpha, \beta) \frac{n_{\mathcal{A} \cup \mathcal{B}}(\alpha) n_{\mathcal{A} \cup \mathcal{B}}(\beta)}{n_{\mathcal{A} \cup \mathcal{B}}} = \end{aligned}$$

$$\begin{aligned}
 &= \sum_{\alpha, \beta \in \Lambda} \rho(\alpha, \beta) \left(\frac{n_A(\alpha)n_A(\beta)}{n_A} + \frac{n_B(\alpha)n_B(\beta)}{n_B} - \right. \\
 &\quad \left. - \frac{(n_A(\alpha) + n_B(\alpha))(n_A(\beta) + n_B(\beta))}{n_A + n_B} \right) = \\
 &= \sum_{\alpha, \beta \in \Lambda} \rho(\alpha, \beta) \frac{(n_A(\alpha)n_B - n_B(\alpha)n_A)(n_A(\beta)n_B - n_B(\beta)n_A)}{n_A n_B (n_A + n_B)}.
 \end{aligned}$$

Если $R \in \mathbb{R}^{L \times L}$ — матрица расстояний: $R_{i,j} = \rho(\lambda_i, \lambda_j)$, и компонентами вектора $t_{A,B} \in \mathbb{R}^L$ являются числа $t_{A,B,i} = n_A(\lambda_i)n_B - n_B(\lambda_i)n_A$, то

$$d(A) + d(B) - d(A \cup B) = \frac{1}{n_A n_B (n_A + n_B)} \langle R t_{A,B}, t_{A,B} \rangle. \quad (7)$$

Заметим, что для любых A и B

$$\sum_{i=1}^L t_{A,B,i} = \sum_{i=1}^L n_A(\lambda_i)n_B - \sum_{i=1}^L n_B(\lambda_i)n_A = n_A n_B - n_B n_A = 0. \quad (8)$$

Этот факт будет использован для доказательства следующего утверждения:

Утверждение 1. Если $\rho(\alpha, \beta)$ — дискретная метрика (3), тогда

$$d(A) + d(B) \leq d(A \cup B) \quad \forall A, B,$$

$d(A) + d(B) = d(A \cup B)$ тогда и только тогда, когда

$$n_A(\lambda_i) = n_B(\lambda_i) \quad \forall \lambda_i \in \Lambda.$$

Доказательство. Легко видеть, что матрица расстояний R для дискретной метрики имеет два различных собственных значения: $(L - 1)$ с собственным вектором $(-1, -1, \dots, -1)$ и -1 степени $L - 1$ с собственным подпространством

$$T \subset \mathbb{R}^L : \sum_{i=1}^L t_i = 0 \quad \forall t \in T.$$

Таким образом, для любых A и B , соответствующий им вектор $t_{A,B}$ всегда будет содержаться в T и

$$\begin{aligned}
 d(A) + d(B) - d(A \cup B) &= \frac{1}{n_A n_B (n_A + n_B)} \langle R t_{A,B}, t_{A,B} \rangle = \\
 &= \frac{-1}{n_A n_B (n_A + n_B)} \|t_{A,B}\|^2. \quad \square
 \end{aligned}$$

Таким образом, объединение кластеров увеличивает неоднородность состава в случае дискретной метрики.

4. Мера неоднозначности классифицирования

Определение 2. Для данного набора $\tilde{\mathcal{A}}$ непересекающихся подмножеств строк A и столбца j , мера неоднозначности классифицирования, $a(\tilde{\mathcal{A}}, j)$, есть

$$a(\tilde{\mathcal{A}}, j) = \sum_{\lambda \in \Lambda} \sum_{A \in \tilde{\mathcal{A}}} \frac{n_A(\lambda, j)(n_{\tilde{\mathcal{A}}}(\lambda, j) - n_A(\lambda, j))}{n_{\tilde{\mathcal{A}}}(\lambda, j)}. \quad (9)$$

Мера неоднозначности классифицирования $a(\tilde{\mathcal{A}}, j)$ может быть проинтерпретирована в терминах теории вероятности следующим образом. Рассмотрим двух игроков. Первый случайным образом выбирает строку s из $\tilde{\mathcal{A}}$ и сообщает символ λ в j -й позиции этой строки другому игроку, при этом храня номер кластера $\tilde{\mathcal{A}}_1$, из которого эта строка была взята, в секрете. Второй игрок, используя информацию о λ , пытается угадать этот кластер. Вероятность того, что второй игрок даст неверный ответ будет равна

$$p(\tilde{\mathcal{A}}_1 \neq \tilde{\mathcal{A}}_2) = \frac{a(\tilde{\mathcal{A}}, j)}{N_{\tilde{\mathcal{A}}}}, \quad (10)$$

где $\tilde{\mathcal{A}}_2$ — предположение 2-го игрока и $N_{\tilde{\mathcal{A}}}$ — общее число строк в наборе $\tilde{\mathcal{A}}$. Таким образом, $a(\tilde{\mathcal{A}}, j)$ пропорционально вероятности неверного классифицирования строки s .

Минимум $a(\tilde{\mathcal{A}}, j) = 0$ достигается, если множества символов на j -й позиции у разных кластеров не пересекаются. Если, напротив, наборы символов в j -той колонке у всех кластеров совпадают,

$$\forall A \in \tilde{\mathcal{A}} \quad n_A(\lambda, j) = \begin{cases} 1, & \lambda \in \{\lambda_1, \lambda_2, \dots, \lambda_l\}, \\ 0, & \text{иначе,} \end{cases}$$

то неоднозначность классифицирования будет максимальной,

$$a(\tilde{\mathcal{A}}, j) = (n_{\tilde{\mathcal{A}}} - 1)l = N_{\tilde{\mathcal{A}}} - l.$$

Если набор $\tilde{\mathcal{A}}'$ получен из набора $\tilde{\mathcal{A}}$ объединением подмножеств строк A и $B \in \tilde{\mathcal{A}}$: $C = A \cup B$, $C \in \tilde{\mathcal{A}}'$, разность в неоднозначности классифицирования этих двух наборов:

$$a(\tilde{\mathcal{A}}') - a(\tilde{\mathcal{A}}) = \sum_{\lambda \in \Lambda} \left(\frac{n_C(\lambda)(n_{\tilde{\mathcal{A}}}(\lambda) - n_C(\lambda))}{n_{\tilde{\mathcal{A}}'}(\lambda)} - \frac{n_A(\lambda)(n_{\tilde{\mathcal{A}}}(\lambda) - n_A(\lambda))}{n_{\tilde{\mathcal{A}}}(\lambda)} - \frac{n_B(\lambda)(n_{\tilde{\mathcal{A}}}(\lambda) - n_B(\lambda))}{n_{\tilde{\mathcal{A}}}(\lambda)} \right) =$$

$$= \sum_{\lambda \in \Lambda} \frac{1}{n_{\tilde{A}}(\lambda)} ((n_{\tilde{A}}(\lambda) - n_A(\lambda) - n_B(\lambda))(n_A(\lambda) + n_B(\lambda)) - n_A(\lambda)(n_{\tilde{A}}(\lambda) - n_A(\lambda)) - n_B(\lambda)(n_{\tilde{A}}(\lambda) - n_B(\lambda))) = - \sum_{\lambda \in \Lambda} \frac{2n_A(\lambda)n_B(\lambda)}{n_{\tilde{A}}(\lambda)} \leq 0.$$

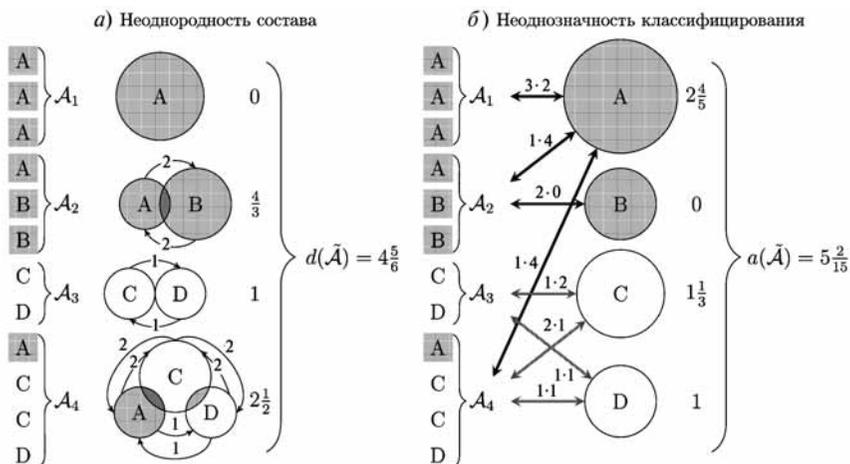
Таким образом, объединение кластеров уменьшает неоднозначность классифицирования.

5. Свертка векторного критерия качества кластеризации

Итак, для любого разбиения \tilde{A} строк матрицы $A \in \Lambda^{M \times N}$ были определены две его характеристики, неоднородность состава и неоднозначность классифицирования:

$$D(\tilde{A}) = \sum_{j=1}^N d(\tilde{A}, j), \quad A(\tilde{A}) = \sum_{j=1}^N a(\tilde{A}, j). \tag{11}$$

Вместе они образуют двухкомпонентный критерий качества кластеризации $Q(\tilde{A}) = (D(\tilde{A}), A(\tilde{A}))$ (см. пример на рис. 1). Чем меньше $D(\tilde{A})$, тем более похожи друг на друга строки в пределах одного кластера. Чем меньше $A(\tilde{A})$, тем сильнее отличается один кластер от другого. Естественно считать кластеризацию оптимальной, если она обладает обоими этими свойствами. Однако, для произвольной $A \in \Lambda^{N \times M}$ мы, вообще говоря,



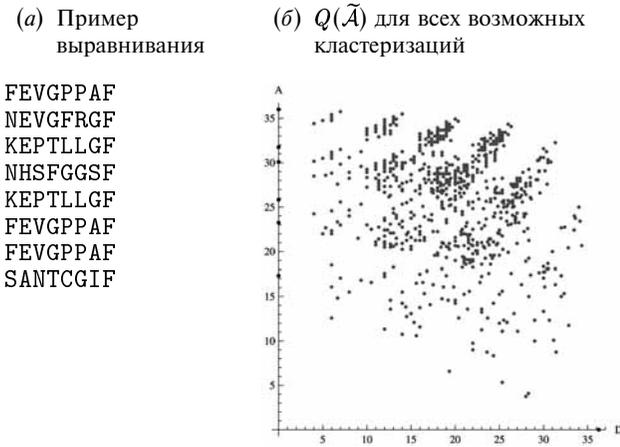


Рис. 2. Пример вычисления $Q(\tilde{A})$

можем не найти такого \tilde{A}^* , которое бы минимизировало одновременно $D(\tilde{A}^*)$ и $A(\tilde{A}^*)$. Например, подразбиение кластеров на более мелкие части уменьшает $D(\tilde{A})$, но увеличивает $A(\tilde{A})$. С другой стороны, объединение кластеров уменьшает $A(\tilde{A})$ и увеличивает $D(\tilde{A})$ (см. рис. 2).

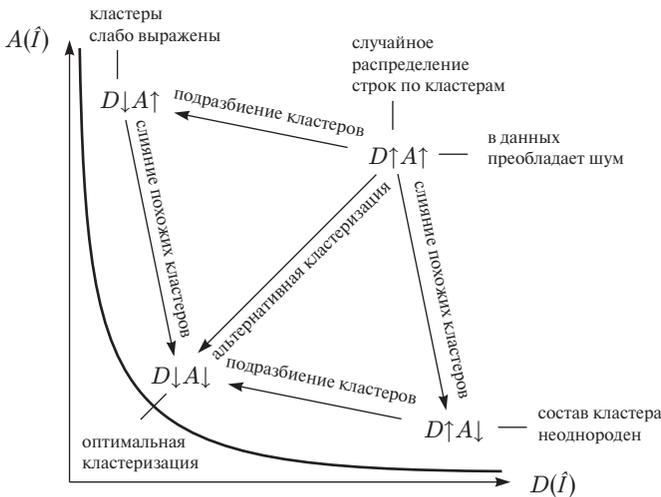


Рис. 3. Множество Парето для $Q(\tilde{A})$ и диаграмма преобразования кластеризаций

Все кластеризации, такие что ни $D(\tilde{A})$, ни $A(\tilde{A})$ не могут быть уменьшены без увеличения другой меры, являются *Парето оптимальными* кластеризациями ($\text{PO}(\mathfrak{A})$) (см. рис. 3).

Для выбора *единственной* оптимальной кластеризации необходимы дополнительные соображения. Одним из возможных решений является рассмотрение некоторой *свертки* функций $D(\tilde{A})$ и $A(\tilde{A})$. Та оптимальная по Парето кластеризация $\tilde{A}^* \in \text{PO}(\mathfrak{A})$, которая будет доставлять минимум этому функционалу, будет считаться оптимальной. Функция

$$q_w(\tilde{A}) = E^2(\tilde{A}) + wC^2(\tilde{A}) \quad w > 0 \quad (12)$$

является одним из примеров такой свертки. По определению, $D(\tilde{A}) \leq M$ и $A(\tilde{A}) \leq M$, поэтому можно положить вес w равным 1. Точка минимума

$$\tilde{A}^* = \operatorname{argmin}_{\tilde{A} \in \text{PO}(\mathfrak{A})} q_1(\tilde{A}) = \operatorname{argmin}_{\tilde{A} \in \mathfrak{A}} q_1(\tilde{A})$$

будет являться проекцией начала координат пространства $D - A$ на образ оптимальных по Парето кластеризаций.

Литература

1. Журавлев Ю., Рязанов В., Сеньков О. «Распознавание». Математические методы. Программная система. Практические применения. М.: Фазис, 2006.
2. Cluster analysis and display of genome-wide expression patterns / M. B. Eisen, P. T. Spellman, P. O. Brown, D. Botstein // Proceedings of the National Academy of Sciences. 1998. V. 95. P. 14863–14868.
3. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs / S. F. Altschul, T. L. Madden, A. A. Schäffer et al. // Nucleic Acids Research. 1997. V. 25. № 17. P. 3389–3402.
4. Henikoff S., Henikoff J. G. Amino acid substitution matrices from protein blocks // Proceedings of the National Academy of Sciences of the United States of America. 1992. V. 89. P. 10915–10919.
5. Reva B., Antipin Y., Sander C. Determinants of protein function revealed by combinatorial entropy optimization. // Genome Biology. 2007. V. 8. № 11. P. R232 (<http://genomebiology.com/2007/8/11/R232>).
6. Thompson J. D., Higgins D. G., Gibson T. J. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice // Nucleic Acids Research. 1994. V. 22. № 22. P. 4673–4680.