

О взаимосвязи мер кластеров и распределении расстояний в компактных метрических пространствах

Пушняков А. С.

Аннотация

Рассматривается компактное метрическое пространство с ограниченной борелевской мерой. Под r -кластером понимается любое измеримое множество диаметра не более r . Исследуется существование набора фиксированного числа $2r$ -кластеров обладающего следующими свойствами: кластеры попарно отделены друг от друга на расстояние r , и мера набора — суммарная мера кластеров набора, — близка к мере всего пространства. Показано, что среди таких наборов существует набор максимальной меры. Для распределения расстояний вводится r -параметрическая дискретизация на *короткие*, *средние* и *длинные* расстояния. В терминах данной дискретизации получена нижняя оценка на меру набора максимальной меры.

Ключевые слова: кластеризация, компактное метрическое пространство, борелевская мера, метрика Хаусдорфа, теорема Бляшке, максимальное паросочетание.

1 Введение

Во многих задачах интеллектуального анализа данных для описания объектов используется метрическая информация [1, 2]. Для задач классификации и кластеризации [3, 4, 5] предполагается, что используемая метрика удовлетворяет так называемому *принципу компактности*: близкие объекты скорее должны лежать в одном классе, нежели в разных [6, 7]. В случае *хорошей* метрики можно полагать, что множество объектов распадается на несколько кластеров, отделенных друг от друга.

Для метрик, представимых в виде объединения кластеров, распределение расстояний имеет некоторые характерные особенности. Если у метрики есть некоторое характерное внутрикластерное расстояние r и межкластерное расстояние R , причем $R > r$, то доля расстояний в промежутке (r, R) мала. В данной статье исследуется следующий вопрос: что нужно потребовать от распределения расстояний, чтобы гарантировать наличие кластерной структуры в метрике?

Мы рассматриваем компактные пространства с ограниченной борелевской мерой (или, что тоже самое, компактные метрические тройки Громова [8, 9]). В данных

терминах удобно определить кластерную структуру как набор фиксированного числа отделенных друг от друга кластеров. Ниже будет показано, что среди таких структур найдется структура максимальной меры. Тогда факт *близости* меры данной структуры к мере всего пространства можно интерпретировать как представление метрики в виде объединения кластеров.

Следуя идее внутрикластерных и межкластерных расстояний, все расстояния классифицируются на *короткие, средние* и *длинные*, и мы требуем, чтобы доля *средних* расстояний была мала. Также мы потребуем некоторые дополнительные ограничения на распределение расстояний, обусловленные количеством кластеров в искомой кластерной структуре. В терминах параметрических ограничений на распределение расстояний мы получим нижнюю оценку на меру кластерной структуры максимальной меры. Вначале мы докажем искомую оценку для конечных полуметрических пространств с равномерной мерой, а затем, используя теорему Бляшке [10], обобщим оценку на случай компактного пространства.

2 Постановка задачи

Пусть дано компактное метрическое пространство (X, ρ) и ограниченная борелевская мера μ на X . Любое борелевское подмножество X диаметра не более r будем называть r -кластером.

Определение 1. Семейство $2r$ -кластеров $\mathcal{X} = \{X_1, \dots, X_k\}$ будем называть r -кластерной структурой порядка k , если $\rho(X_i, X_j) \geq r$ при всех $1 \leq i < j \leq k$, где $\rho(A, B) = \inf\{\rho(x, y) : x \in A, y \in B\}$. Мерой \mathcal{X} назовем величину $\mu(\mathcal{X}) \stackrel{\text{def}}{=} \sum_{i=1}^k \mu(X_i)$.

Верно следующее

Утверждение 1. Среди всех r -кластерных структур порядка k есть структура максимальной меры.

Доказательство. Достаточно рассматривать структуры, содержащие только компактные множества. Для каждой структуры произвольно занумеруем входящие в неё множества. Пусть

$$\mu^* = \sup_{\mathcal{X}} \mu(\mathcal{X}) < \infty$$

Тогда найдется последовательность $\{\mathcal{X}_n\}_{n=1}^{\infty}$ такая, что $\lim_{n \rightarrow \infty} \mu(\mathcal{X}_n) = \mu^*$. По теореме Бляшке метрическое пространство компактов из X по метрике Хаусдорфа ρ_H является компактным. Тогда без ограничения общности можно считать $X_{in} \xrightarrow{\rho_H} X_i^*$, где $\mathcal{X}_n = \{X_{in}\}_{i=1}^k$. Очевидно, что $\mathcal{X}^* = \{X_i^*\}_{i=1}^k$ является r -кластерной структурой порядка k . Пусть $\mu^* - \mu(\mathcal{X}) = \delta > 0$. Рассмотрим множества $X_i^\varepsilon = \{y : \rho(y, X_i^*) \leq \varepsilon\}$. При достаточно малых ε

$$\sum_{i=1}^k (\mu(X_i^\varepsilon) - \mu(X_i^*)) \leq \frac{\delta}{2}, \quad \mu^* - \sum_{i=1}^k \mu(X_i^\varepsilon) > \frac{\delta}{2}.$$

Но в силу сходимости $X_{in} \xrightarrow{\rho_H} X_i^*$ при фиксированном ε и достаточно больших n выполнено $X_{in} \subset X_i^\varepsilon$, и

$$\mu(\mathcal{X}_n) = \sum_{i=1}^k \mu(X_{in}) \leq \sum_{i=1}^k \mu(X_i^\varepsilon) < \mu^* - \frac{\delta}{2}.$$

Получили противоречие. \square

Нашей дальнейшей задачей является определение условий, при которых можно гарантировать, что отношение $\frac{\mu(\mathcal{X}^*)}{\mu(X)}$ близко к единице, где \mathcal{X}^* — r -кластерная структура порядка k максимальной меры. Везде далее мы считаем $k \geq 2$.

Рассмотрим модельный пример: пусть $X = X_1 \sqcup \dots \sqcup X_k$, где все множества X_i являются r -кластерами, и $\rho(X_i, X_j) \geq R > r$ при всех $1 \leq i < j \leq k$. В данном случае мера $\frac{r}{2}$ -кластерной структуры порядка k максимальной меры равна $\mu(X)$ и выполнены равенства:

$$\{(x_1, \dots, x_{k+1}) \in X^{k+1} : \rho(x_i, x_j) > r, 1 \leq i < j \leq k+1\} = \emptyset$$

$$\{(x, y) \in X^2 : r < \rho(x, y) < R\} = \emptyset$$

Пару точек $(x, y) \in X^2$ будем называть *ребром*, длина ребра — это $\rho(x, y)$ (более наглядная аналогия будет видна в случае $|X| < \infty$). В выше описанном примере нет ребер длины которых лежат в интервале (r, R) , а также среди любых $k+1$ точек есть ребро длины не больше r .

Если $\rho(x, y) \leq r$, то будем называть ребро (x, y) *r -коротким*; если $\rho(x, y) > 3r$, то будем называть ребро (x, y) *r -длинным*; все остальные ребра — *r -средние*. Набор точек (x_1, \dots, x_k) назовем *r -антикликкой порядка k* , если $\rho(x_i, x_j) > r$ при всех $1 \leq i < j \leq k$. Если понятно, о каком r идет речь, то приставка r будет опускаться.

В нижеследующих неравенствах (1) и (2) пары (x, y) и наборы (x_1, \dots, x_{k+1}) мы считаем упорядоченными. Потребуем, чтобы в нашем метрическом пространстве (X, ρ) мера r -средних ребер была *мала* в следующем смысле:

$$M(X) = \frac{1}{2} \mu\{(x, y) \in X^2 : r < \rho(x, y) \leq 3r\} \leq \frac{1}{2} \alpha \mu(X)^2, \quad (1)$$

а мера r -антиклик порядка $k+1$ была *мала* в следующем смысле:

$$T_{k+1}(X) = \frac{1}{(k+1)!} \mu\{(x_1, \dots, x_{k+1}) \in X^{k+1} : \rho(x_i, x_j) > r, 1 \leq i < j \leq k+1\} \leq \frac{1}{(k+1)!} \beta \mu(X)^{k+1}, \quad (2)$$

где $\alpha, \beta > 0$ — параметры (мы будем их считать достаточно малыми).

Далее мы докажем, что при выполнении условия (1) и (2) верна оценка меры \mathcal{X}^* вида

$$\mu(\mathcal{X}^*) \geq \Psi(\alpha, \beta) \mu(X), \quad (3)$$

где $\Psi(\alpha, \beta) \rightarrow 1$ при $\alpha \rightarrow 0$ и $\beta \rightarrow 0$.

Выбор верхней границы для интервала средних ребер объясняется техническими соображениями, и, гипотетически, может быть уменьшен. Нижняя же граница увеличена быть не может, если мы хотим получить оценку вида (3). Это следует из следующего утверждения.

Утверждение 2. Для любых $\delta > 0$ и $\alpha > 0$ существует компактное метрическое пространство (X, ρ) такое, что

$$\mu\{(x, y) \in X^2 : r < r' < \rho(x, y)\} \leq \alpha \mu(X)^2,$$

и мера любого $2r$ кластера не более $\delta \mu(X)$.

Доказательство. Пусть $X = Y_1 \sqcup \dots \sqcup Y_s$, и $|Y_i| = m$ при $1 \leq i \leq s$, а мера любого множества равна его мощности. Определим расстояние ρ следующим образом:

$$\rho(x, y) = \begin{cases} 2r', & x, y \in Y_i \\ r', & x \in Y_i, y \in Y_j, i \neq j \end{cases}$$

Тогда мера любого $2r$ -кластера не превосходит s , а

$$\mu\{(x, y) \in X^2 : r < r' < \rho(x, y)\} = sm(m - 1)$$

Осталось взять $m > \frac{1}{\delta}$, $s > \frac{1}{\alpha}$. □

3 Жадная кластерная структура

Вначале мы получим оценку вида (3) в случае, когда (X, ρ) — конечное полуметрическое пространство. Отметим, что достаточно получить нижнюю оценку меры какой-то r -кластерной структуры порядка k . Рассмотрим следующую жадную процедуру. Пусть X_1 — множество максимальной мощности среди всех $2r$ -кластеров (если таких множеств несколько, то выберем любое). Обозначим его окрестность r за Z_1 , т.е.

$$Z_1 = \{x \in X : \rho(x, X_1) < r\}$$

Пусть у нас есть попарно непересекающиеся множества Z_1, \dots, Z_m . Тогда X_{m+1} — множество максимальной мощности среди всех $2r$ -кластеров в $X \setminus \bigcup_{i=1}^m Z_i$, а множество

Z_{m+1} — r -окрестность X_{m+1} во множестве $X \setminus \bigcup_{i=1}^m Z_i$, т.е.

$$Z_{m+1} = \left\{ x \in X \setminus \bigcup_{i=1}^m Z_i : \rho(x, X_{m+1}) < r \right\}$$

Так как мощность X конечна, то процедура оборвётся на некотором шаге.

Определение 2. Построенное разбиение $X = \bigsqcup_{i=1}^n Z_i$ мы назовем жадным кластерным разбиением, а семейство $2r$ -кластеров $\{X_1, \dots, X_k\}$ назовем жадной r -кластерной структурой порядка k .

Сделаем несколько замечаний относительно последнего определения. Во-первых, последовательности Z_i и X_i определяются неоднозначно — далее считается, что фиксирована некоторая пара последовательностей (X_i, Z_i) . Во-вторых, из построения очевидно, что жадная r -кластерная структура порядка k является r -кластерной структурой порядка k по определению (1).

Отметим, что последовательность $\{|X|_i\}_{i=1}^n$ монотонно убывает, однако, для последовательности $\{|Z|_i\}_{i=1}^n$ свойство монотонности в общем случае не выполняется. Пусть $\{W_i\}_{i=1}^n$ — упорядоченные по убыванию $|Z|_i$. Следующим шагом мы покажем, что в условиях (1) и (2) и при достаточно малых α и β первые k по мощности Z_i покрывают почти все множество X , т.е. верно неравенство

$$\sum_{i=1}^k W_i \geq \Phi(\alpha, \beta)|X|, \quad (4)$$

где $\Phi(\alpha, \beta) \rightarrow 1$ при $\alpha \rightarrow 0$ и $\beta \rightarrow 0$.

4 Нижняя оценка числа антиклик

Пусть $T_s(i_1, \dots, i_s)$ — число r -антиклик порядка s , таких, что ровно по одной вершине содержится в каждом из множеств Z_{i_j} . Понятно, что при $s = 1$ выполнено $T_s(i_j) = |Z_{i_j}|$. Нам понадобится следующее рекуррентное соотношение на $T_s(i_1, \dots, i_s)$.

Утверждение 3. Пусть $i_1 < \dots < i_s$, то при $s \geq 2$

$$T_s(i_1, \dots, i_s) \geq \frac{|Z_{i_1}|}{s} T_{s-1}(i_2, \dots, i_s), \quad (5)$$

Доказательство. Пусть x_2, \dots, x_s — вершины некоторой антиклики, $x_j \in Z_{i_j}$. Для каждой из вершин x_j рассмотрим множества

$$S(x_j) = \{y \in Z_{i_1} : \rho(x_j, y) \leq r\},$$

Так как диаметр $S(x_j)$ не более $2r$, то $|S(x_j)| \leq |Z_{i_1}|$. Пусть $Y = Z_{i_1} \setminus \bigcup_{j=2}^s S(x_j)$, тогда

$$|Y| \geq \frac{|Z_{i_1}|}{s}.$$

Для любой точки $y \in Y$ вершины y, x_2, \dots, x_s образуют r -антиклику порядка s . Тогда имеем

$$T_s(i_1, \dots, i_s) \geq \sum_{(x_2, \dots, x_s)} \frac{|Z_{i_1}|}{s} = \frac{|Z_{i_1}|}{s} T_{s-1}(i_2, \dots, i_s)$$

□

Из утверждения (3) сразу же получаем

$$T_s(i_1, \dots, i_s) \geq \frac{1}{s!} \prod_{j=1}^s |Z_{i_j}| \quad (6)$$

Теперь мы получим нижнюю оценку на $T_{k+1}(X)$ — число r -антиклик порядка $k+1$. Нам осталось только просуммировать неравенство (6) по всем наборам из $k+1$ множеств Z_i . Введем обозначение для симметрического многочлена от n переменных

$$\sigma_s(y_1, \dots, y_n) \stackrel{\text{def}}{=} \sum_{1 \leq i_1 < \dots < i_s \leq n} \prod_{j=1}^s y_{i_j}, \quad (7)$$

тогда, используя (2) и (6), получим

$$\frac{1}{(k+1)!} \sigma_{k+1}(z_1, \dots, z_n) \leq T_{k+1}(X) \leq \frac{1}{(k+1)!} \beta |X|^{k+1} \quad (8)$$

Разделим каждое z_i на $|X|$ и упорядочим по убыванию: получим набор $w_1 \geq \dots \geq w_n$, и тогда

$$\sigma_{k+1}(w_1, \dots, w_n) \leq \beta \quad (9)$$

5 Нижняя оценка для $\sum_{j=1}^k W_j$

По сути мы получили следующую задачу оптимизации

$$\left\{ \begin{array}{l} f(\mathbf{w}) = \sum_{j=1}^k w_j \rightarrow \min_{\mathbf{w}} \\ w_i \geq 0 \\ w_i \geq w_j, \quad i \leq j \\ \sum_{i=1}^n w_i = 1 \\ \sigma_{k+1}(w_1, \dots, w_n) \leq c \end{array} \right. \quad (10)$$

Очевидно, что задача (10) имеет решение. Мы будем далее считать, что $n > k$, иначе решение задачи (10) очевидно. Нам понадобятся следующие простые утверждения.

Утверждение 4. Пусть \mathbf{w} — решение задачи (10), тогда либо $w_i = w_j$ при всех $1 \leq i \leq j \leq n$, либо $\sigma_{k+1}(w_1, \dots, w_n) = c$.

Доказательство. Предположим противное. Пусть $w_l > w_{l+1}$. Рассмотрим вектор

$$\mathbf{w}^\varepsilon = (w_1^\varepsilon, \dots, w_n^\varepsilon) = (w_1 - (n-l)\varepsilon, \dots, w_l - (n-l)\varepsilon, w_{l+1} + l\varepsilon, \dots, w_n + l\varepsilon)$$

Так как $\sigma_{k+1}(w_1, \dots, w_n) < c$, то при достаточно малых $\varepsilon > 0$ вектор \mathbf{w}^ε будет допустимым для задачи (10). Но $f(\mathbf{w}^\varepsilon) < f(\mathbf{w})$, получили противоречие. \square

Утверждение 5. Пусть \mathbf{w} — решение задачи (10) и $w_k = \lambda > 0$. Тогда $\mathbf{w} = (w_1, \underbrace{\lambda, \dots, \lambda}_s, \mu, 0, \dots, 0)$, где $s \geq k-1$ и $\mu < \lambda$.

Доказательство. Если все w_i попарно равны λ , то утверждение верно. Для любых $1 \leq i < j \leq n$ имеем

$$\begin{aligned} \sigma_{k+1}(w_1, \dots, w_n) &= w_i w_j \sigma_{k-1}(w_1, \dots, \hat{w}_i, \dots, \hat{w}_j, \dots, w_n) + \\ &+ (w_j + w_i) \sigma_k(w_1, \dots, \hat{w}_i, \dots, \hat{w}_j, \dots, w_n) + \sigma_{k+1}(w_1, \dots, \hat{w}_i, \dots, \hat{w}_j, \dots, w_n) \end{aligned}$$

Пусть нашлось $2 \leq i \leq k$ такое, что $w_i > \lambda = w_{i+1} = \dots = w_k$. Тогда рассмотрим вектор

$$\mathbf{w}' = (w_1 + w_i - \lambda, w_2, \dots, w_{i-1}, \lambda, \dots, \lambda)$$

$$\sigma_{k+1}(\mathbf{w}) - \sigma_{k+1}(\mathbf{w}') = (w_1 w_i - (w_1 + w_i - \lambda) \lambda) \sigma_{k-1}(w_1, \dots, \hat{w}_i, \dots, \hat{w}_j, \dots, w_n) > 0$$

Пусть w_l — последняя ненулевая компонента \mathbf{w} . Пусть j — первая компонента \mathbf{w} такая, что $w_j < \lambda$, и $j < l$. Тогда рассмотрим вектор

$$\mathbf{w}' = (w_1, \dots, w_{j-1}, \min\{\lambda, w_l + w_j\}, w_{j+1}, \dots, w_{l-1}, \max\{w_l + w_j - \lambda, 0\}, \dots, 0)$$

$$\sigma_{k+1}(\mathbf{w}) - \sigma_{k+1}(\mathbf{w}') = (w_j w_l - \min\{\lambda, w_l + w_j\} \max\{w_l + w_j - \lambda, 0\}) \sigma_{k-1}(w_1, \dots, \hat{w}_i, \dots, \hat{w}_j, \dots, w_n) > 0$$

Тогда по утверждению (4) получаем, что w — не решение задачи (10). \square

Итак, рассмотрим вектор $\mathbf{w} = (w_1, \underbrace{\lambda, \dots, \lambda}_s, \mu, 0, \dots, 0)$.

$$\begin{aligned} \sigma_{k+1}(\mathbf{w}) &= w_1 \mu \sigma_{k-1}(\underbrace{\lambda, \dots, \lambda}_s) + (w_1 + \mu) \sigma_k(\underbrace{\lambda, \dots, \lambda}_s) + \sigma_{k+1}(\underbrace{\lambda, \dots, \lambda}_s) = \\ &= w_1 \mu \binom{s}{k-1} \lambda^{k-1} + (w_1 + \mu) \binom{s}{k} \lambda^k + \binom{s}{k+1} \lambda^{k+1} \leq c \end{aligned} \quad (11)$$

Мы рассмотрим несколько случаев.

1. $s = k - 1$. Неравенство (11) переходит в

$$\frac{1}{k+1}\mu\lambda^{k-1} \leq w_1\mu\lambda^{k-1} \leq c$$

$$f(\mathbf{w}) = 1 - \mu \geq 1 - \min\left\{\lambda, \frac{c(k+1)}{\lambda^{k-1}}\right\} \geq 1 - (c(k+1))^{\frac{1}{k}}$$

2. $s = k$.

$$\frac{1}{k+1}\lambda^k \leq w_1\mu k\lambda^{k-1} + (w_1 + \mu)\lambda^k \leq c, \quad \lambda \leq (c(k+1))^{\frac{1}{k}}$$

$$f(\mathbf{w}) \geq 1 - 2\lambda \geq 1 - 2(c(k+1))^{\frac{1}{k}}$$

3. $s \geq k + 1$

$$\left(\frac{s\lambda}{k+1}\right)^{k+1} \leq \binom{s}{k+1}\lambda^{k+1} \leq c, \quad s\lambda \leq (k+1)c^{\frac{1}{k+1}}$$

$$f(\mathbf{w}) \geq 1 - s\lambda \geq 1 - (k+1)c^{\frac{1}{k+1}}$$

Так как при $k \geq 2$ и $c \leq \frac{1}{2}$ выполнено $(k+1)c^{\frac{1}{k+1}} \geq 2(c(k+1))^{\frac{1}{k}}$, то верно следующее

Утверждение 6. Пусть \mathbf{w} — решение задачи (10) и $(k+1)c^{\frac{1}{k+1}} \leq 1$. Тогда

$$f(\mathbf{w}) \geq 1 - (k+1)c^{\frac{1}{k+1}}. \quad (12)$$

Используя соотношения (9) и (12) получаем

$$\sum_{i=1}^k W_i \geq |X| \left(1 - (k+1)\beta^{\frac{1}{k+1}}\right) \stackrel{\text{def}}{=} \Phi(\alpha, \beta)|X|, \quad (13)$$

6 Оценка меры жадной r -кластерной структуры

Далее мы будем рассматривать только внутреннюю структуру множеств Z_i . Поэтому без ограничения общности можно полагать, что $|Z_i| = W_i$.

Нам осталось доказать, что

$$\sum_{i=1}^k |Z_i| - \sum_{i=1}^k |X_i| = o(1)|X|, \quad \alpha + \beta \rightarrow 0$$

Рассмотрим множества Z_i и $X_i \subset Z_i$. Для любых $x \in X_i$ и $z \in Z_i$ выполнено $\rho(x, z) \leq 3r$, поэтому концы всех длинных ребер лежат в $Z_i \setminus X_i$. Рассмотрим во множестве $Z_i \setminus X_i$ максимальное паросочетание из длинных ребер, которое покрывает множество W_i . Пусть $Y_i = Z_i \setminus (X_i \cup W_i)$, тогда $X_i \cap Y_i$ является $3r$ -кластером. Докажем простое утверждение, связывающее мощность X_i и число средних ребер в $X_i \cap Y$.

Утверждение 7. Пусть (A, ρ) — конечное полуметрическое пространство диаметра не более $3r$, и множество B является $2r$ -кластером максимальной мощности. Тогда число средних ребер не менее $M(A) \geq \frac{1}{2}|A||A \setminus B|$.

Доказательство. Пусть x_0 — точка, из которой выходит максимально число коротких ребер, а S — замкнутый шар радиуса r с центром в x_0 . Тогда

$$M(A) \geq \frac{1}{2}|A|(|A| - |S|) \geq \frac{1}{2}|A|(|A| - |B|)$$

□

Также для любого ребра (u_1, u_2) из паросочетания, покрывающего W_i , и точки $x \in X_i$ хотя бы одно из ребер (x, u_j) является средним. В купе с утверждением (7) получаем следующее неравенство:

$$M(Z_i) \geq \frac{1}{2}(|X_i| + |Y_i|)|Y_i| + \frac{1}{2}|W_i||X_i| \quad (14)$$

Сейчас мы применим технику аналогичную той, что использовалась при оценке числа антиклик.

Утверждение 8. Пусть $T_s(Z_i)$ — число r -антиклик порядка s во множестве Z_i . Тогда при $s \geq 3$

$$T_s(Z_i) \geq \frac{1}{s}(|Z_i| - (s-1)|X_i|)_+ T_{s-1}(Z_i)$$

Доказательство. Доказательство почти дословно совпадает с доказательством утверждения (3). Пусть x_1, \dots, x_{s-1} образуют некоторую антиклику. Для каждой из вершин x_j рассмотрим множества

$$S(x_j) = \{y \in Z_i : \rho(x_j, y) \leq r\},$$

Так как диаметр $S(x_j)$ не более $2r$, то $|S(x_j)| \leq |X_{i_1}|$. Пусть $Y = Z_i \setminus \bigcup_{j=2}^s S(x_j)$, тогда

$$|Y| \geq (|Z_i| - (s-1)|X_i|)_+.$$

Для любой точки $y \in Y$ вершины y, x_1, \dots, x_{s-1} образуют антиклику порядка s . Осталось заметить что каждую антиклику порядка s мы посчитали не более s раз, тогда имеем

$$T_s(Z_i) \geq \sum_{(x_1, \dots, x_{s-1})} \frac{1}{s}(|Z_i| - (s-1)|X_i|)_+ = \frac{1}{s}(|Z_i| - (s-1)|X_i|)_+ T_{s-1}(Z_i).$$

□

Из утверждения (6) и равенства $T_1(Z_i) = |Z_i|$ сразу следует неравенство

$$T_{k+1}(Z_i) \geq \frac{1}{(k+1)!} \prod_{j=1}^{k+1} (|Z_i| - (j-1)|X_i|)_+ \quad (15)$$

Если $|X_i|(k+1) \leq |Z_i|$, то

$$T_{k+1}(Z_i) \geq \left(\frac{|Z_i|}{k+1} \right)^{k+1}$$

Пусть I_1 — множество всех индексов от $1 \leq i \leq k$ таких, что $|X_i|(k+1) \leq |Z_i|$, тогда

$$\begin{aligned} \left(\frac{\sum_{i \in I_1} |Z_i|}{k(k+1)} \right)^{k+1} &\leq \frac{1}{(k+1)^{k+1}} \sum_{i \in I_1} |Z_i|^{k+1} \leq \sum_{i \in I_1} T_{k+1}(Z_i) \leq T_{k+1}(X) \leq \frac{\beta |X|^{k+1}}{(k+1)!} \\ &\sum_{i \in I_1} |Z_i| \leq ek\beta^{\frac{1}{k+1}} |X| \end{aligned}$$

Если же $|X_i|(k+1) > |Z_i|$, то из неравенства (14) получаем

$$\begin{aligned} M(Z_i) &\geq \frac{1}{2(k+1)} |Z_i| (|Y_i| + |W_i|) = \frac{1}{2(k+1)} |Z_i| (|Z_i| - |X_i|) \\ |Z_i| - |X_i| &\leq \frac{2(k+1)M(Z_i)}{|Z_i|} \end{aligned}$$

Рассмотрим I_2 — множество таких индексов от $1 \leq i \leq k$, что $i \notin I_1$ и $|Z_i| \geq \sqrt{\alpha} |X|$. Тогда суммируя предыдущее неравенство по множеству I_2 :

$$\sum_{i \in I_2} (|Z_i| - |X_i|) \leq \frac{2(k+1)M(X)}{\sqrt{\alpha}|X|} \leq \sqrt{\alpha}(k+1)|X|$$

Наконец, получаем

$$\begin{aligned} \sum_{i=1}^k (|Z_i| - |X_i|) &\leq \sum_{i \in I_1} |Z_i| + \sum_{i \in I_2} (|Z_i| - |X_i|) + \sum_{i \notin I_1 \cup I_2} |Z_i| \leq \\ &\leq ek\beta^{\frac{1}{k+1}} |X| + \sqrt{\alpha}(k+1)|X| + \sqrt{\alpha}k|X| = (\sqrt{\alpha}(2k+1) + ke\beta^{\frac{1}{k+1}}) |X| \end{aligned}$$

Итак, мы доказали следующую теорему

Теорема 1. Пусть (X, ρ) конечное полуметрическое пространство с равномерной мерой μ , а \mathcal{X}^* — r -кластерная структура максимальной меры. Тогда, если выполнены неравенства (1) и (2), то

$$\mu(\mathcal{X}^*) \geq \Psi(\alpha, \beta) |X|, \quad (16)$$

где

$$\Psi(\alpha, \beta) = 1 - \sqrt{\alpha}(2k+1) - (k(e+1) + 1)\beta^{\frac{1}{k+1}}$$

7 Обобщение на случай произвольного компактного пространства

Мы будем использовать технику, аналогичную той, что была использована при доказательстве утверждения (1).

Теорема 2. Пусть (X, ρ) компактное метрическое пространство с ограниченной борелевской мерой μ , а \mathcal{X}^* — r -кластерная структура максимальной меры. Тогда, если выполнены неравенства (1) и (2), то выполнено неравенство (16).

Доказательство. Фиксируем произвольное $0 < \varepsilon < 1$. В X существует конечная ε -сеть, а значит и разбиение X на конечное число N_ε ε -кластеров $\{A_i\}_{i=1}^{N_\varepsilon}$. Выберем N_ε положительных рациональных чисел $q_1, \dots, q_{N_\varepsilon}$ так, что $\mu(A_i) \geq q_i$ при $1 \leq i \leq N_\varepsilon$ и $q_i \geq \mu(A_i)(1 - \varepsilon)$.

Рассмотрим полуметрическое пространство конечной мощности $X_\varepsilon = B_1 \sqcup \dots \sqcup B_s$, где $\frac{|B_i|}{|B_j|} = \frac{q_i}{q_j}$, а функция расстояния ρ_ε определяется следующим образом:

$$\rho_\varepsilon(x, y) = \begin{cases} 0, & x, y \in B_i \\ \rho(A_i, A_j), & x \in B_i, y \in B_j, i \neq j \end{cases}$$

Отметим, что

$$|B_i| = \frac{q_i |X_\varepsilon|}{\sum_{j=1}^{N_\varepsilon} q_j} \leq \frac{\mu(A_i) |X_\varepsilon|}{(1 - \varepsilon) \sum_{j=1}^{N_\varepsilon} \mu(A_j)} = \frac{\mu(A_i) |X_\varepsilon|}{(1 - \varepsilon) \mu(X)}$$

Если $x \in B_i, y \in B_j$ и $\rho_\varepsilon(x, y) \in (r, 3r]$, то для всех $v \in A_i, u \in A_j$ верно $\rho(v, u) \in (r, 3r + 2\varepsilon]$. Отсюда получаем оценку на число r -средних ребер в X_ε :

$$\begin{aligned} M(X_\varepsilon) &= \sum_{1 \leq i < j \leq N_\varepsilon} [\rho(A_i, A_j) > r] |B_i| |B_j| \leq \\ &\leq \frac{|X_\varepsilon|^2}{(1 - \varepsilon)^2 \mu(X)^2} \sum_{1 \leq i < j \leq N_\varepsilon} [\rho(A_i, A_j) > r] \mu(A_i) \mu(A_j) \leq \\ &\leq \frac{|X_\varepsilon|^2}{2(1 - \varepsilon)^2 \mu(X)^2} \mu\{(x, y) \in X^2 : r < \rho(x, y) \leq 3r + 2\varepsilon\} \stackrel{(1)}{\leq} \\ &\leq \frac{|X_\varepsilon|^2}{2(1 - \varepsilon)^2} \left(\alpha + \frac{1}{\mu(X)^2} \mu\{(x, y) \in X^2 : 3r < \rho(x, y) \leq 3r + 2\varepsilon\} \right) \stackrel{\text{def}}{=} \frac{1}{2} \alpha_\varepsilon |X_\varepsilon|^2 \end{aligned}$$

Аналогично имеем оценку для r -антиклик порядка $k + 1$:

$$T_{k+1}(X_\varepsilon) = \sum_{1 \leq i_1 < \dots < i_{k+1} \leq N_\varepsilon} \prod_{1 \leq j < l \leq k+1} [\rho(A_{i_j}, A_{i_l}) > r] \prod_{j=1}^{k+1} |B_{i_j}| \leq$$

$$\begin{aligned} &\leq \frac{|X_\varepsilon|^{k+1}}{(1-\varepsilon)^{k+1}\mu(X)^{k+1}} \sum_{1 \leq i_1 < \dots < i_{k+1} \leq N_\varepsilon} \prod_{1 \leq j < l \leq k+1} [\rho(A_{i_j}, A_{i_l}) > r] \prod_{j=1}^{k+1} \mu(A_{i_j}) \stackrel{(2)}{\leq} \\ &\leq \frac{|X_\varepsilon|^{k+1} \beta}{(k+1)!(1-\varepsilon)^{k+1}} \stackrel{\text{def}}{=} \frac{1}{(k+1)!} \beta_\varepsilon |X_\varepsilon|^{k+1} \end{aligned}$$

Заметим, что при $\varepsilon \rightarrow 0$ $\alpha_\varepsilon \rightarrow \alpha$ и $\beta_\varepsilon \rightarrow \beta$. В силу теоремы 1 получаем, что в X_ε существует r -кластерная структура порядка k $\mathcal{C}_\varepsilon = \{C_{1\varepsilon}, \dots, C_{k\varepsilon}\}$ меры не менее $\Psi(\alpha_\varepsilon, \beta_\varepsilon)|X_\varepsilon|$.

Понятно, что каждое B_i либо полностью содержится в каком-то множестве семейства \mathcal{C}_ε , либо никакой элемент B_i не входит ни в какое множество семейства \mathcal{C}_ε . Для каждого $C_{i\varepsilon} = B_{j_1} \sqcup \dots \sqcup B_{j_l}$ рассмотрим множество $X_{i\varepsilon} = cl(A_{j_1} \sqcup \dots \sqcup A_{j_l})$ в X . Заметим, что множество $X_{i\varepsilon}$ является $(r + 2\varepsilon)$ -кластером, и для любых $1 \leq i < j \leq k$ выполнено $\rho(X_{i\varepsilon}, X_{j\varepsilon}) \geq r$.

Настало время снова применить теорему Бляшке. Пусть $\varepsilon \rightarrow 0$. Рассматривая последовательность наборов $(X_{1\varepsilon}, \dots, X_{k\varepsilon})$, без ограничения общности можно считать, что $X_{i\varepsilon} \xrightarrow{\rho_H} X_i^*$. Очевидно, что множества X_i^* образуют r -кластерную структуру порядка k , назовем её \mathcal{X}^* . Более того

$$\sum_{i=1}^k \mu(X_{i\varepsilon}) \geq \frac{(1-\varepsilon)\mu(X)}{|X_\varepsilon|} \sum_{i=1}^k |C_{i\varepsilon}| \geq (1-\varepsilon)\Psi(\alpha_\varepsilon, \beta_\varepsilon)\mu(X) \xrightarrow{\varepsilon \rightarrow 0} \Psi(\alpha, \beta)\mu(X).$$

Почти дословно повторяя доказательство утверждения (1), получаем $\mu(\mathcal{X}^*) \geq \Psi(\alpha, \beta)\mu(X)$. \square

Список литературы

- [1] Журавлев Ю. И., Никифоров В. В. Алгоритмы распознавания, основанные на вычислении оценок // *Кибернетика*. — 1971. — no. 3. — Pp. 1–11.
- [2] Айзерман М. А., Браверман Э. М., Розоноэр Л. И. Метод потенциальных функций в теории обучения машин. — Наука, 1970.
- [3] Celebi M. E., Kingravi H. A., Vela P. A. A comparative study of efficient initialization methods for the k-means clustering algorithm // *Expert Systems with Applications*. — 2013. — Vol. 40, no. 1. — Pp. 200–210.
- [4] De Amorim R. C., Mirkin B. Minkowski metric, feature weighting and anomalous cluster initializing in k-means clustering // *Pattern Recognition*. — 2012. — Vol. 45, no. 3. — Pp. 1061–1075.
- [5] Aggarwal C. C., Reddy C. K. Data clustering: algorithms and applications. — CRC Press, 2013.

- [6] *Загоруйко Н. Г.* Гипотезы компактности и λ -компактности в методах анализа данных // *Сибирский журнал индустриальной математики*. — 1998. — Vol. 1, no. 1. — Pp. 114–126.
- [7] *Браверман Э. М.* Опыты по обучению машины распознаванию зрительных образов // *Автоматика и телемеханика*. — 1962. — Vol. 23, no. 3. — Pp. 349–365.
- [8] *Gromov M.* Metric structures for Riemannian and non-Riemannian spaces. — Springer Science & Business Media, 2007.
- [9] *Вершик А. М.* Универсальное пространство урысона, метрические тройки громова и случайные метрики на натуральном ряде // *Успехи математических наук*. — 1998. — Vol. 53, no. 5 (323). — Pp. 57–64.
- [10] *Половинкин Е. С., Балашов М. В.* Элементы выпуклого и сильно выпуклого анализа. — Физматлит, 2004.