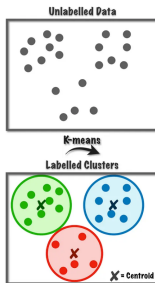


# Оценка качества кластеризации

Виктор Китов

[victorkitov.github.io](https://victorkitov.github.io)

Курс поддержан  
фондом  
'Интеллект'



Победитель  
конкурса VK среди  
курсов по IT



# Оценка качества кластеризации

## Оценка качества кластеризации:

- если кластеризация-промежуточный этап: по качеству итоговой задачи
- если нет разметки:
  - используют идею, что кластеризация хороша, если:
    - объекты одного кластера похожи
    - объекты разных кластеров непохожи
- если есть разметка:
  - учитывать инвариантность к переименованию
  - имеет смысл для малого #размеченных объектов
    - иначе - классификация

# Содержание

- 1 Оценки не использующие разметку
- 2 Оценки, использующие разметку

## Метрики качества<sup>1</sup>

- Пусть  $z_n$  - номер кластера для  $x_n$ .
- Среднее внутрикластерное расстояние:

$$F_0 = \frac{\sum_{i < j} \mathbb{I}[z_i = z_j] \rho(x_i, x_j)}{\sum_{i < j} \mathbb{I}[z_i = z_j]}$$

- Среднее межкластерное расстояние:

$$F_1 = \frac{\sum_{i < j} \mathbb{I}[z_i \neq z_j] \rho(x_i, x_j)}{\sum_{i < j} \mathbb{I}[z_i \neq z_j]}$$

- Композитные метрики:

$$F_0/F_1, F_1 - F_0$$

---

<sup>1</sup>Какие метрики нужно максимизировать, а какие - минимизировать?

## Индекс Дэвиса-Болдуина

- $s_i = \frac{1}{|C_i|} \sum_{n \in C_i} \rho(\mu_i, x_n)$  - радиус кластера  $i$ .
- $d_{ij} = \rho(\mu_i, \mu_j)$  - расстояние между центроидами  $i$  и  $j$ .
- Качество разделения кластеров  $i$  и  $j$ :

$$R_{ij} = \frac{s_i + s_j}{d_{ij}}$$

- Индекс Дэвиса-Болдуина:

$$DB = \frac{1}{K} \sum_{k=1}^K \max_{i \neq k} R_{ik}$$

- ⊕ : Быстро вычисляется.
- ⊖ : Поощряет выпуклые кластера
- ⊖ : тип расстояния определяет  $\mu_i$

## Коэффициент силуэта<sup>2</sup>

Качество кластеризации каждого объекта  $x_i$ :

$$Silhouette_i = \frac{d_i - s_i}{\max\{d_i, s_i\}}$$

где среднее расстояние от  $x_i$  до объектов

- $s_i$  - того же кластера
- $d_i$  - ближайшего чужого кластера

Общее качество классификации (коэффициент силуэта):

$$Silhouette = \frac{1}{N} \sum_{i=1}^N \frac{d_i - s_i}{\max\{d_i, s_i\}}$$

---

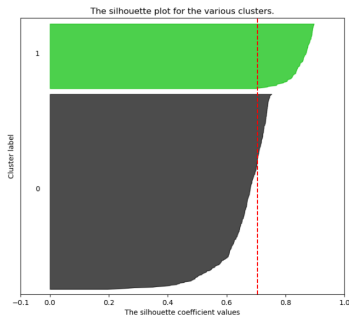
<sup>2</sup>Rousseeuw (1987). "Silhouettes: a Graphical Aid to the Interpretation and Validation of Cluster Analysis". Computational and Applied Mathematics 20: 53–65.

# Обсуждение

- Преимущества
  - Интерпретируемость:  $Silhouette \in [-1, 1]$ ,
    - 1: идеальная кластеризация
    - 0: случайная кластеризация
    - -1: полностью некорректная (инвертированная) кластеризация
- Недостатки
  - сложность  $O(N^2 D)$ 
    - можно рассчитывать по случайной подвыборке
  - поощряет выпуклые кластеры

## Подбор #кластеров по силуэту<sup>3</sup>

- Отсортируем объекты в каждом кластере по коэффициенту силуэта.
- Качество кластеризации - среднее значение коэффициента и отсутствие отрицательных значений.



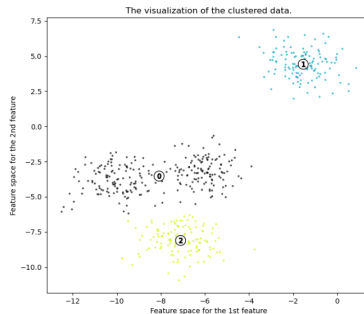
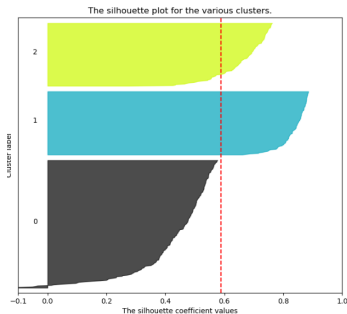
---

<sup>3</sup>Эксперимент в sklearn.



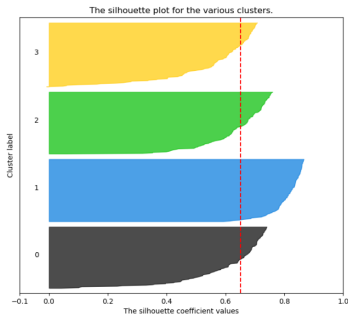
## Подбор #кластеров по силуэту<sup>3</sup>

- Отсортируем объекты в каждом кластере по коэффициенту силуэта.
- Качество кластеризации - среднее значение коэффициента и отсутствие отрицательных значений.



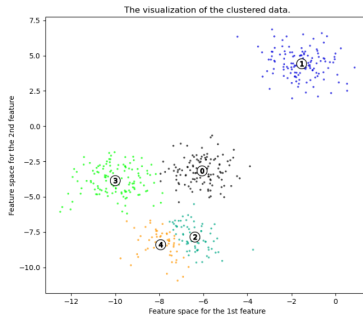
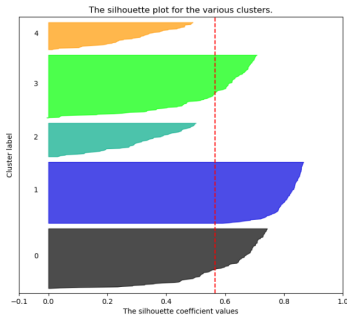
## Подбор #кластеров по силуэту<sup>3</sup>

- Отсортируем объекты в каждом кластере по коэффициенту силуэта.
- Качество кластеризации - среднее значение коэффициента и отсутствие отрицательных значений.



## Подбор #кластеров по силуэту<sup>3</sup>

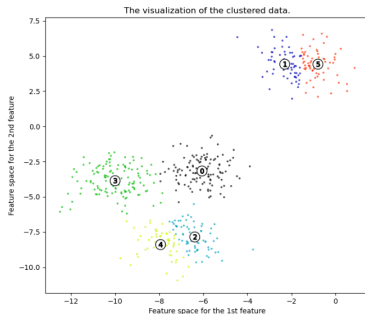
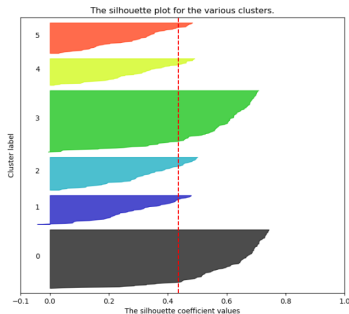
- Отсортируем объекты в каждом кластере по коэффициенту силуэта.
- Качество кластеризации - среднее значение коэффициента и отсутствие отрицательных значений.



<sup>3</sup>Эксперимент в sklearn.

## Подбор #кластеров по силуэту<sup>3</sup>

- Отсортируем объекты в каждом кластере по коэффициенту силуэта.
- Качество кластеризации - среднее значение коэффициента и отсутствие отрицательных значений.



## Индекс Калинского<sup>4</sup>

- Внутрикластерная (within cluster) ковариационная матрица

$$W = \frac{1}{N - K} \sum_{k=1}^K \sum_{x \in C_k} (x - \mu_k) (x - \mu_k)^T$$

- Межкластерная (between cluster) ковариационная матрица

$$B = \frac{1}{K - 1} \sum_{k=1}^K N_k (\mu_k - \mu) (\mu_k - \mu)^T$$

- Индекс Калинского:

$$I = \frac{\text{tr } B}{\text{tr } W} = \frac{N - K}{K - 1} \frac{\text{tr} \left\{ \sum_{k=1}^K N_k (\mu_k - \mu) (\mu_k - \mu)^T \right\}}{\text{tr} \left\{ \sum_{k=1}^K \sum_{x \in C_k} (x - \mu_k) (x - \mu_k)^T \right\}}$$

---

<sup>4</sup>[https://www.researchgate.net/publication/233096619\\_A\\_Dendrite\\_Method\\_for](https://www.researchgate.net/publication/233096619_A_Dendrite_Method_for)

## Индекс Калинского

- Используем свойства

$$\sum_i \operatorname{tr} \{ \alpha_i A_i \} = \sum_i \alpha_i \operatorname{tr} A_i$$

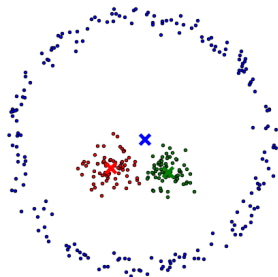
$$\operatorname{tr} \{ AB \} = \operatorname{tr} \{ BA \}, \quad \operatorname{tr} a = a \quad \forall a \in \mathbb{R}$$

$$\begin{aligned} I &= \frac{\operatorname{tr} B}{\operatorname{tr} W} = \frac{N - K}{K - 1} \frac{\operatorname{tr} \left\{ \sum_{k=1}^K N_k (\mu_k - \mu) (\mu_k - \mu)^T \right\}}{\operatorname{tr} \left\{ \sum_{k=1}^K \sum_{x \in C_k} (x - \mu_k) (x - \mu_k)^T \right\}} \\ &= \frac{N - K}{K - 1} \frac{\sum_{k=1}^K N_k \operatorname{tr} \left\{ (\mu_k - \mu)^T (\mu_k - \mu) \right\}}{\sum_{k=1}^K \sum_{x \in C_k} \operatorname{tr} \left\{ (x - \mu_k) (x - \mu_k)^T \right\}} = \frac{N - K}{K - 1} \frac{\sum_{k=1}^K N_k \|\mu_k - \mu\|^2}{\sum_{k=1}^K \sum_{x \in C_k} \|x - \mu_k\|^2} \end{aligned}$$

- Измеряем отношение межкластерного к внутрикластерному разбросу.

## Ограничение для невыпуклого кластера

- Сложность  $O(ND)$ , но поощряет выпуклые кластеры.
- Здесь качество Калинского будет казаться низким, как и индекс Дэвида-Болдуина:



- $\sum_{k=1}^K N_k \|\mu_k - \mu\|^2$  мало, а  $\sum_{k=1}^K \sum_{x \in C_k} \|x - \mu_k\|^2$  велико
- Коэффициент силуэта будет вести себя лучше.

## Алгоритм Monti consensus clustering<sup>5</sup>

- Генерируем  $H$  псевдовыборок  $D_1, D_2, \dots, D_H$  из  $X$ , кластеризуем каждую.
- На  $D_h$   $(x_i, x_j)$  кластеризуются как  $(z_i^h, z_j^h)$ ,  $h \in M(i, j)$ .
  - $M(i, j) = \{h : x_i \in D_h \ \& \ x_j \in D_h\}$
- Определим матрицу консенсуса (consensus matrix)

$$M(i, j) = \frac{\sum_{h \in M(i, j)} \mathbb{I}[z_i^h = z_j^h]}{|M(i, j)|}, \quad M \in \mathbb{R}^{N \times N}$$

- $M(i, j) \in \{0, 1\} \Rightarrow$  у точек  $(i, j)$  устойчивая класт-ция.
- $M(i, j) \in (0, 1) \Rightarrow$  у точек  $(i, j)$  неустойчивая класт-ция.
  - тем неустойчивее, чем ближе  $M(i, j)$  к 0.5.

<sup>5</sup>[https://en.wikipedia.org/wiki/Consensus\\_clustering](https://en.wikipedia.org/wiki/Consensus_clustering)

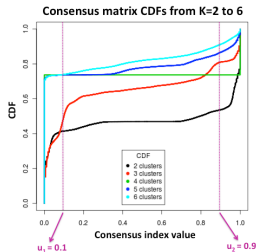


# Алгоритм Monti consensus clustering

- Для каждого #кластеров  $K$  посчитаем пропорцию пар точек с неопределённой кластеризацией (proportion of ambiguous clustering, PAC)

$$PAC(K) = \frac{|\{(i, j) : i < j \text{ \& } 0.1 \leq M(i, j) \leq 0.9\}|}{C_2^N}$$

- $PAC(K) \in [0, 1]$  - мера неустойчивости кластеризации на  $K$  кластеров;  $K^* = \arg \min_K PAC(K)$ .



PAC for each  $K$  is  $CDF_K(u_2) - CDF_K(u_1)$ .  
 According to this criterion, optimal  $K$  here is 4.

# Содержание

- 1 Оценки не использующие разметку
- 2 Оценки, использующие разметку

## Перекрестная таблица

- Пример перекрестной таблицы (contingency matrix):

	кластер 1	кластер 2	кластер 3
класс 1	5	2	0
класс 2	0	3	4

- $\oplus$  : Определяем разброс каждого класса по кластерам и разброс кластера по классам.
- $\ominus$  : Сложно анализировать для большого числа кластеров/классов. Не числовая метрика качества.

## Перекрестная таблица

- Пример перекрестной таблицы (contingency matrix):

	кластер 1	кластер 2	кластер 3
класс 1	5	2	0
класс 2	0	3	4

- ⊕ : Определяем разброс каждого класса по кластерам и разброс кластера по классам.
- ⊖ : Сложно анализировать для большого числа кластеров/классов. Не числовая метрика качества.
- Числовая мера качества - Unsupervised Clustering Accuracy:
  - $\Pi$  - всевозможные перенумеровки номеров кластеров

$$ACC(c, z) = \max_{\pi \in \Pi} \frac{1}{N} \sum_{n=1}^N \mathbb{I}[c_n = \pi(z_n)]$$

## Матрица сочетаемости

- Матрица сочетаемости  $\in \mathbb{R}^{2 \times 2}$  вычисляет счётчики  $\#$  пар  $(x_i, x_j)$ .

	$z_i = z_j$	$z_i \neq z_j$
$y_i = y_j$	$n_{11}$	$n_{12}$
$y_i \neq y_j$	$n_{21}$	$n_{22}$

- Как понять по матрице качество кластеризации?

## Матрица сочетаемости

- Матрица сочетаемости  $\in \mathbb{R}^{2 \times 2}$  вычисляет счётчики  $\# \text{пар}$   $(x_i, x_j)$ .

	$z_i = z_j$	$z_i \neq z_j$
$y_i = y_j$	$n_{11}$	$n_{12}$
$y_i \neq y_j$	$n_{21}$	$n_{22}$

- Как понять по матрице качество кластеризации?
- $\oplus$  : Определяем сочетаемость разбиения по классам-кластерам.
- $\ominus$  : Не числовая метрика качества. Какую предложим?

## Rand index

- Rand index - единая метрика по матрице сочетаемости.
- Пусть  $y_1, \dots, y_N$  - истинная разметка. Обозначим<sup>6</sup>

$$n_{11} = |\{(x_i, x_j) : z_i = z_j \ \& \ y_i = y_j\}|$$

$$n_{22} = |\{(x_i, x_j) : z_i \neq z_j \ \& \ y_i \neq y_j\}|$$

$$\text{RandInd} = RI = \frac{n_{11} + n_{22}}{C_2^N} = \frac{n_{11} + n_{22}}{n_{11} + n_{12} + n_{21} + n_{22}} \in [0, 1]$$

- В чем недостаток?

---

<sup>6</sup> Это loss или score?

<sup>7</sup> [Adjusted Rand Index - wikipedia.](#)

<sup>8</sup> J-близость Жаккарда между множеством пар, у которых совпали классы и множеством пар, у которых совпали кластеры.

## Rand index

- Rand index - единая метрика по матрице сочетаемости.
- Пусть  $y_1, \dots, y_N$  - истинная разметка. Обозначим<sup>6</sup>

$$n_{11} = |\{(x_i, x_j) : z_i = z_j \ \& \ y_i = y_j\}|$$

$$n_{22} = |\{(x_i, x_j) : z_i \neq z_j \ \& \ y_i \neq y_j\}|$$

$$\text{RandInd} = RI = \frac{n_{11} + n_{22}}{C_2^N} = \frac{n_{11} + n_{22}}{n_{11} + n_{12} + n_{21} + n_{22}} \in [0, 1]$$

- В чем недостаток?  $\uparrow RI$  с  $\uparrow \#$ кластеров. Лучше<sup>7,8</sup>

$$\text{AdjRandInd} = \frac{RI - \mathbb{E}\{RI\}}{\max(RI) - \mathbb{E}\{RI\}} \text{ или } J = \frac{n_{11}}{n_{11} + n_{12} + n_{21}}$$

<sup>6</sup> Это loss или score?

<sup>7</sup> [Adjusted Rand Index - wikipedia.](#)

<sup>8</sup> J-близость Жаккарда между множеством пар, у которых совпали классы и множеством пар, у которых совпали кластеры.



## Гомогенность<sup>9</sup>

- Обозначим  $N = \# \text{объектов}$ ,  $n_k = \# \text{объектов в кластере } k$ ,  $m_c = \# \text{объектов в классе } c$ ,  $n_{ck} = \# \text{объектов класса } c \text{ в кластере } k$ .

$$H_{class} = - \sum_{c=1}^C \frac{m_c}{N} \log \frac{m_c}{N}$$

$$H_{clust} = - \sum_{k=1}^K \frac{n_k}{N} \log \frac{n_k}{N}$$

$$H_{class|clust} = - \sum_{k=1}^K \frac{n_k}{N} \sum_{c=1}^C \frac{n_{ck}}{n_k} \log \frac{n_{ck}}{n_k}$$

$H_{class|clust} = 0$  при полном объяснении,  $H_{class|clust} = 1$  нет связи

---

<sup>9</sup><https://aclanthology.org/D07-1043.pdf>

## Гомогенность

$$\text{Homogeneity} = 1 - \frac{H(class|clust)}{H(class)}$$

- Гомогенность показывает долю информации о классах, объясненной кластеризацией.
  - 1: в кластерах представители только 1 класса
  - 0: в кластерах распределение классов=априорному распределению
- Какой недостаток?

## Гомогенность

$$\text{Homogeneity} = 1 - \frac{H(class|clust)}{H(class)}$$

- Гомогенность показывает долю информации о классах, объясненной кластеризацией.
  - 1: в кластерах представители только 1 класса
  - 0: в кластерах распределение классов=априорному распределению
- Какой недостаток? Гомогенность поощряет  $\uparrow \#$  кластеров
  - $=1$ , когда каждый объект - в своём кластере

Полнота<sup>10</sup>

- Нужна доп. мера полноты (насколько объекты одного класса оказываются в одном кластере)

$$\text{Completeness} = 1 - \frac{H(\text{clust}|\text{class})}{H(\text{clust})}$$

- Полнота = 1, если класс полностью определяет кластер (все объекты кластера-в одном классе)
- Какой недостаток?

---

<sup>10</sup><https://aclanthology.org/D07-1043.pdf>

Полнота<sup>10</sup>

- Нужна доп. мера полноты (насколько объекты одного класса оказываются в одном кластере)

$$\text{Completeness} = 1 - \frac{H(\text{clust}|\text{class})}{H(\text{clust})}$$

- Полнота =1, если класс полностью определяет кластер (все объекты кластера-в одном классе)
- Какой недостаток? Полнота поощряет ↓#кластеров
  - =1, когда все объекты в одном кластере

---

<sup>10</sup><https://aclanthology.org/D07-1043.pdf>

V-мера<sup>11</sup>

- V-мера - среднее гармоническое от гомогенности и полноты.

$$V = \frac{1}{\frac{1}{2} \frac{1}{\text{Homogeneity}} + \frac{1}{2} \frac{1}{\text{Completeness}}}$$

- Взвешенный учёт гомогенности и полноты:

$$V_{\beta} = \frac{1}{\left(\frac{\beta}{1+\beta}\right) \frac{1}{\text{Homogeneity}} + \frac{1}{1+\beta} \frac{1}{\text{Completeness}}}$$

- $V = V_{\beta}$  при  $\beta = 1$ .

---

<sup>11</sup><https://aclanthology.org/D07-1043.pdf>

## Нормализованная взаимная информация

- Взаимная информация - степень связи сл. вел.  $X, Y$ :

$$\begin{aligned}
 MI(X, Y) &= KL(P(X, Y) || P(X)P(Y)) \\
 &= \sum_{x \in \text{dom}(X)} \sum_{y \in \text{dom}(Y)} P(x, y) \log \frac{P(x, y)}{P(x)P(y)}
 \end{aligned}$$

$$MI(X, Y) = H(Y) - H(Y|X) = H(X) - H(X|Y)$$

- Нормализованная взаимная информация ( $NMI \in [0, 1]$ )<sup>12</sup>  
- др. вариант агрегации полноты и гомогенности:

$$\begin{aligned}
 NMI(clust, class) &= \frac{MI(clust, class)}{\max\{H_{clust}, H_{class}\}} \\
 &= \frac{H(clust) - H(clust|class)}{\max\{H_{clust}, H_{class}\}} = \frac{H(class) - H(class|clust)}{\max\{H_{clust}, H_{class}\}}
 \end{aligned}$$

<sup>12</sup>Это loss или score?

## Заключение

- Оценки, не использующие разметку:
  - размеры кластеров vs. межкластерные расстояния
- Оценки, использующие разметку:
  - сопоставление кластеров с истинными метками
    - важна инвариантность к перенумеровке кластеров