Функции расстояния. Меры сходства

- Функция расстояния
- Виды функции расстояния
- Матрица сходства
- Меры сходства

Для того чтобы решить задачу кластерного анализа, построить алгоритм классификации, необходимо количественно определить функции расстояния, понятия сходства (близости) и разнородности (отдаленности) объектов, кластеров между собой.

Что означает утверждение " два объекта S_i , S_j различны"?

Это означает, что если объекты S_i , S_j попадают в один и тот же класс (при построении кластеров различными алгоритмами), когда расстояние (отдаленность) между соответствующими объектами было бы "достаточно малым" в определенном смысле, и, наоборот, попадают в разные классы, если расстояние между S_i , S_j будет "достаточно большим". Таким образом, определяем понятие расстояния между объектами S_i , S_j из исходного множества допустимых объектов S.

Если объекты S_i S являются точками n - мерного пространства, то исходное множество объектов S будет подмножеством евклидового пространства E_p .

Определение 2.1. Неотрицательная вещественнозначная функция $\rho(S_i, S_j)$ называется функцией расстояния (метрикой), если выполнимы аксиомы:

- 1. (S_i, S_i) 0 для всех S_i, S_i из S_i
- 2. $(S_i, S_j)=0$ тогда и только тогда, когда $S_i=S_j$;
- 3. $(S_i, S_i) = (S_i, S_i)$;
- 4. (S_i, S_i) (S_i, S_k) + (S_k, S_i) , где S_i , S_i и S_k любые вектора из E_v .

Значение (S_i, S_i) для заданных S_i, S_i называется расстоянием между объектами S_i, S_i .

Аксиомы 1,2 определяют положительность метрики, аксиома 3 - определяет свойство симметричности, аксиома 4 - неравенство треугольника.

Евклидова метрика очень популярна и наиболее употребительна. Если S_i =($_{i1}$, $_{i2}$,..., $_{in}$) и S_j =($_{j1}$, $_{j2}$,..., $_{jn}$), то евклидова метрика представляется как

$$\rho(S_i, S_i) = \mathbf{\dot{c}}. \tag{2.1}$$

При использовании этого расстояния будем считать, что компоненты (признаки) S_i взаимно независимы, однородны по своему физическому смыслу и все они одинаково важны с точки зрения вопроса об отнесении объекта к тому или иному классу, а также признаковое пространство совпадает с геометрическим пространством и понятие близости объектов соответственно совпадает с понятием геометрической близости в этом пространстве.

"Взвешенное" евклидово расстояние. Обычно применяется в ситуациях, когда некоторые признаки содержательно важны в смысле формирования кластеров, поэтому удается приписать каждому из признаков некоторый неотрицательный "вес" _і, "пропорциональный" степени его важности

$$\rho(S_i, S_j) = \sqrt{\omega_{\square^i} (\alpha_1 - \beta_1)^2 + \omega_2 (\alpha_2 - \beta_2)^2 + \dots + \omega_n (\alpha_n - \beta_n)^2}$$
 (2.2)

Конечно, определение весов i, i=1, n связано с дополнительным исследованием, таким как организация опроса экспертов и обработки их мнений, исследования обучающей выборки, использования специальных процедур предварительной обработки.

Хемингово расстояние. Используется как мера различия объектов, задаваемых дихотомическими признаками. Хемингово расстояние задается с помощью формулы

$$\rho_H(S_i, S_j) = \sum_{K=1}^n |\alpha_{ik} - \beta_{jk}|, \qquad (2.3)$$

откуда видно что оно равно числу $_{ij}$ несовпадений значений соответствующих признаков в рассматриваемых S_{i} -м и S_{i} -м объектах.

Также можно использовать в качестве меры близости объектов S_i и S_j величину $(S_i, S_j) = \frac{1}{i} / n$, где n - число признаков.

Также используются меры типа l-норма, наиболее простая с вычислительной точки зрения

$$\rho\left(S_{i}, S_{j}\right) = \sum_{k=1}^{n} \left|\alpha_{ik} - \beta_{jk}\right| \tag{2.3.1}$$

для произвольных значений признаков.

Сюпремум норма также легко вычисляется, включает себе процедуру упорядочения, и может использоваться в качестве меры близости

$$(2.4)$$

Наиболее часто используемыми (двумерный случай) в задачах классификации изображений и анализа сцен являются метрика абсолютного значения $d_A(S_1,S_2)$ и метрика максимального значения $d_M(S_1,S_2)$:

$$d_{A}(S_{1}, S_{2}) = |\alpha_{1} - \beta_{1}| + |\alpha_{2} - \beta_{2}|,$$

$$d_{M}(S_{1}, S_{2}) = \max \left| \left\langle \Box |\alpha_{1} - \beta_{1}|, |\alpha_{2} - \beta_{2}| \right| \right|.$$

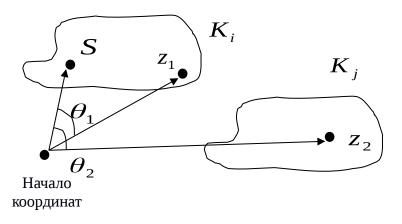
Выбор соответствующей метрики имеет под собой более глубокие основания, чем может показаться на первый взгляд. Это связано с формой представления в ЭВМ информации о графическом изображении и его геометрических параметрах. Так, например, любые точки изображения в цифровом представлении определяются их целочисленными координатами в поле кадра. Результат вычисления евклидова расстояния между парой точек в общем случае не является целым. Если использовать здесь округление, то, как легко убедиться на примере, можно нарушить свойство 4 определения метрики. В то же время, используя метрику абсолютного значения, мы всегда получаем целые расстояния при целочисленных координатах точек. Можно, отметить, что метрика d_A более проста в вычислительном отношении.

Несимметрическая функция сходства основания на косинус угла используется в тех случаях, когда кластеры обнаруживают тенденцию располагаться вдоль главных осей, как это показано на рис. 2.3.

Функция

$$K(S,Z) = \cos \theta = \frac{S'Z}{\|S\|\|Z\|}$$
 (2.5)

представляет собой косинус угла, образованного вектором S и Z, и достигающая максимума, когда их направления совпадают. Этой мерой сходства удобно пользоваться, когда кластеры достаточно отстоят друг от друга и от начала координат.



$$K(S,z_1) = \cos \theta_1 = \frac{S'z_1}{\|S\|\|z_1\|}, K(S,z_2) = \cos \theta_2 = \frac{S'z_2}{\|S\|\|z_2\|}.$$

Из рис. 2.3 видно, что образ z_1 обладает большим сходством с образом S, чем образ z_2 , поскольку значение $K(S,z_1)$ больше значения $K(S,z_2)$.

Когда рассматриваются двоичные образы и их элементы принимают значение из множества [0,1], функции сходства (2.5) можно дать интересную негеометрическую интерпретацию.

Если α_i =1, считается, что двоичный объект S обладает i-м признаком. В таком случае член S'Zв формуле (2.5) просто характеризует число общих для образов S и z признаков, а $\|S\|\|z\|=\sqrt{|S'S|(z'z)}$ - среднее геометрическое число признаков, которыми обладает объект S, и число признаков, которыми обладает объект z. Понятно, что K(S,z) есть мера общих признаков у объектов S и z.

Мера Танимото является двоичным вариантом формулы (2.5), который нашел широкое распространение в информационной классификации болезней и таксономии (классификация видов животных и растений)

$$K(S,z) = \frac{S'z}{S'S + z'z - S'z}$$
 (2.5.1)

Расстояние **Махалонобиса** оказывается полезной мерой сходства в тех случаях, когда статистические характеристики объектов присутствуют в явном виде

$$D = (S - z)'C^{-1}(S - z),$$

где C- коварационная матрица совокупности объектов; z- вектор средних значений; а S - объект.

Данная функция Махалонобиса успешно используется, когда известны плотности распределения объектов в классах, т.е. когда используются статистические свойства классов.

Исходное множество объектов

$$S = S_1, S_2, ..., S_m, S_i = (i_1, i_2, ..., i_n)$$

может быть представлено в виде матрицы исходных данных размером т п:

$$S = [S_{1}, S_{2}, \dots, S_{m}] = \begin{vmatrix} \alpha_{11} & \alpha_{12} & \dots & \alpha_{1n} \\ \alpha_{21} & \alpha_{22} & \dots & \alpha_{2n} \\ \dots & \dots & \dots & \dots \\ \alpha_{m1} & \alpha_{m2} & \dots & \alpha_{mn} \end{vmatrix}.$$

Аналогичным образом расстояния между парами объектов могут быть представлены в виде симметричной матрицы расстояний:

$$C=\mathcal{L}$$
 (2.6)

Если для элементов матрицы расстояний C, применить некоторый порог , т. е.:

$$C_{ij} = \begin{bmatrix} 1 & \square & \square \\ 1 & \text{если } \rho_{ij} < \rho^{\square}, i, j = 1, 2, ..., m \end{bmatrix}$$
 (2.7)

тогда получаем симметричную матрицу сходства $C = \|C_{ii}\|_{m \times m}$, где $C_{ii} = 1$, $C_{ij} = C_{ji}$, $C_{ij} = 0$, 1.

Элемент C_{ij} матрицы сходства C называется мерой сходства объектов S_i , S_j относительно выбранного порога ρ^\square .

Воспользуемся теперь введенным понятием расстояния для вычисления меры рассеяния или разнородности множества объектов

$$S = S_1, S_2, ..., S_m$$
.

Определение 2.2. Пусть задано множество

 $S = S_1, S_2, ..., S_m$ допустимых объектов (множество наблюдений).

Величина $r_{\rho} = \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \rho(S_i, S_j)$ называется общим рассеянием, соответствующим данной функции расстояния $\rho(S_i, S_j)$.

Определение 2.3. Величина $r_{\rho} = \frac{r_{\rho}}{M_{\rho}}$, где $M_{\rho} = \left(\frac{m^2 - m}{2}\right)$ называется средним рассеянием множества S.

Обоснование определений 2.2, 2.3 следует из рассмотрения матрицы расстояний $C = \prod_{ij} | I_{m \times m}$ с учетом того, что, во-первых, для всех i, $\rho_{ii} = \rho(S_i, S_i)$, а во-вторых, из $(S_i, S_j) = (S_i, S_i)$ следует $I_{ij} = I_{ij}$ для всех $I_{ij} = I_{ij$

Отсюда величина r представляет собой сумму m^2 расстояний, из которых m равны нулю, и $\frac{m^2-m}{2}$ различны и неотрицательны. Поэтому r есть арифметическое среднее ненулевых расстояний между парами объектов из S. Также применяется другой вид рассеяния.

Определение 2.4. Матрица $R_3 = \sum_{i=1}^m \left(S_i - \bar{S} \right)$ іі размером n n называется матрицей рассеяния множества S, причем $\bar{S} = \sum_{i=1}^m S_i m$ - есть вектор арифметических средних размером n 1.

Понятие меры внутренней однородности и меры разнородности классов иногда является необходимым условием при построении некоторых алгоритмов кластерного анализа.

Пусть $K_1 = S_1$, S_2 ,..., S_{m1} и $K_2 = S_{m1+1}$, S_{m1+2} ,..., S_{m2} обозначают два класса объектов. Ниже приводятся наиболее употребительные меры близости между классами объектов.

Определение 2.5. Величину

$$D_I(K_I, K_2)$$
=min (S_i, S_j) , для всех $i=1, 2, ..., m_I$, $j=m_I+1, m_I+2, ..., m_2$

будем называть минимальным локальным расстоянием между кластерами K_1 и K_2 , соответствующим данной функции расстояния . Другими словами, расстояние измерено по принципу "ближний сосед" БС.

Определение 2.6. Величину

$$D_2(K_1, K_2) = \max (S_i, S_j)$$
, для всех $i=1, 2, ..., m_1$, $j=m_1+1, m_1+2, ..., m_2$

назовем максимальным локальным расстоянием между K_1 и K_2 или расстоянием, измеренным по принципу " $\partial aльний coced$ " ДС.

Определение 2.7. Величина

$$D_{3}(K_{1},K_{2}) = \sum_{i=1}^{m_{2}} \sum_{j=1}^{m_{1}} \frac{\rho(S_{i},S_{j})}{m_{1} \cdot m_{2}}$$

называется средним расстоянием между K_1 и K_2 относительно $\rho(S_i, S_j)$ или расстоянием измеренным по принципу " $cpe \partial he \ddot{u}$ centure centu