

МАШИННОЕ ОБУЧЕНИЕ

метрические методы классификации, регрессии и кластеризации

Воронцов Константин Вячеславович
ФУПМ МФТИ • ВМК МГУ • Яндекс • FORECSYS

8 июля 2016
Сочи, Сириус • Проектная смена • 1–24 июля 2016

1 Классификация

- Гипотезы непрерывности и компактности
- Обобщённый метрический классификатор
- Связь метрического и линейного классификатора

2 Регрессия

- Метод наименьших квадратов
- Ядерное сглаживание
- Выбор ядра и ширины окна

3 Кластеризация

- Задача кластеризации
- Метод k -средних
- Иерархическая кластеризация

Восстановление зависимости по эмпирическим данным

Задача восстановления зависимости $y = f(x)$
по точкам *обучающей выборки* (x_i, y_i) , $i = 1, \dots, \ell$.

Дано: векторы $x_i = (x_i^1, \dots, x_i^n)$ — объекты обучающей выборки,
 $y_i = f(x_i)$ — правильные ответы, $i = 1, \dots, \ell$:

$$\begin{pmatrix} x_1^1 & \dots & x_1^n \\ \dots & \dots & \dots \\ x_\ell^1 & \dots & x_\ell^n \end{pmatrix} \xrightarrow{f} \begin{pmatrix} y_1 \\ \dots \\ y_\ell \end{pmatrix}$$

Найти: функцию $a(x)$, способную давать правильные ответы
на *тестовых объектах* $\tilde{x}_i = (\tilde{x}_i^1, \dots, \tilde{x}_i^n)$, $i = 1, \dots, q$:

$$\begin{pmatrix} \tilde{x}_1^1 & \dots & \tilde{x}_1^n \\ \dots & \dots & \dots \\ \tilde{x}_q^1 & \dots & \tilde{x}_q^n \end{pmatrix} \xrightarrow{a?} \begin{pmatrix} a(\tilde{x}_1) \\ \dots \\ a(\tilde{x}_q) \end{pmatrix}$$

Гипотезы компактности и непрерывности

Гипотеза непрерывности (для регрессии):

близким объектам соответствуют близкие ответы.

Гипотеза компактности (для классификации):

близкие объекты, как правило, лежат в одном классе.

Формализация понятия «близости»:

задана функция расстояния $\rho(x, x_i)$.

Пример. Евклидово расстояние и его обобщение:

$$\rho(x, x_i) = \left(\sum_{j=1}^n |x^j - x_i^j|^2 \right)^{1/2} \quad \rho(x, x_i) = \left(\sum_{j=1}^n w_j |x^j - x_i^j|^p \right)^{1/p}$$

$x = (x^1, \dots, x^n)$ — вектор признаков объекта x ,

$x_i = (x_i^1, \dots, x_i^n)$ — вектор признаков объекта x_i .

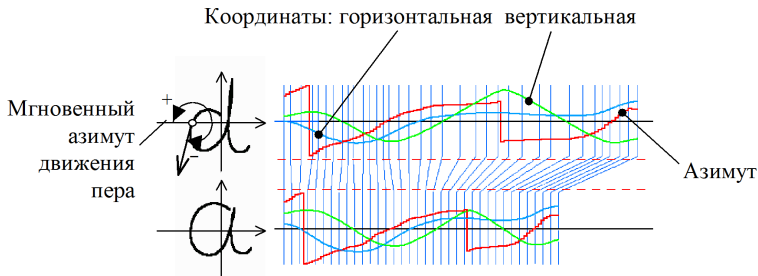
Беспризнаковые способы вычисления расстояний

Расстояния на основе выравниваний:

- между текстами (редакторское расстояние Левенштейна):

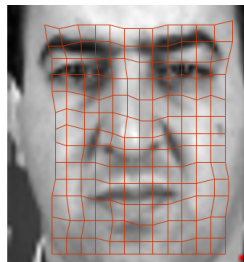
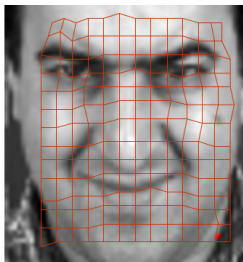
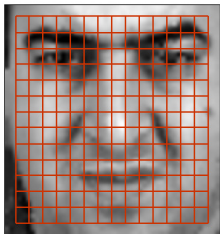
CTGGGCTAAAAGGTCCCTTAGCC--TTTAGAAAAA-GGGCCATTAGGAAAATTGCA
CTGGGACTAAA----CCTTAGCCTATTTACAAAAATGGGCCATTAGG---TTGCA

- между сигналами (энергия сжатий и растяжений):



Беспризнаковые способы вычисления расстояний

Расстояния между изображениями на основе выравнивания:



Оценивается энергия растяжения прямоугольной сетки

Обобщённый метрический классификатор

Для произвольного x отранжируем объекты x_1, \dots, x_ℓ :

$$\rho(x, x^{(1)}) \leq \rho(x, x^{(2)}) \leq \dots \leq \rho(x, x^{(\ell)}),$$

$x^{(i)}$ — i -й сосед объекта x среди x_1, \dots, x_ℓ ;

$y^{(i)}$ — ответ на i -м соседе объекта x .

Метрический алгоритм классификации:

$$a(x; X^\ell) = \arg \max_{y \in Y} \underbrace{\sum_{i=1}^{\ell} [y^{(i)} = y] w(i, x)}_{\Gamma_y(x)},$$

$w(i, x)$ — вес, оценка сходства объекта x с его i -м соседом, неотрицательная, не возрастающая по i .

$\Gamma_y(x)$ — оценка близости объекта x к классу y .

Метод k ближайших соседей (k nearest neighbors, k NN)

$w(i, x) = [i \leq 1]$ — метод ближайшего соседа

$w(i, x) = [i \leq k]$ — метод k ближайших соседей

Преимущества:

- простота реализации (lazy learning);
- параметр k можно оптимизировать по критерию скользящего контроля (leave-one-out):

$$\text{LOO}(k, X^\ell) = \sum_{i=1}^{\ell} [a(x_i; X^\ell \setminus \{x_i\}, k) \neq y_i] \rightarrow \min_k.$$

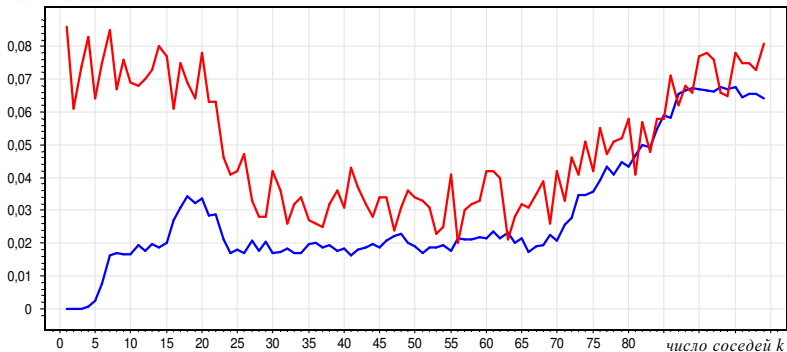
Недостатки:

- неоднозначность классификации при $\Gamma_y(x) = \Gamma_s(x)$, $y \neq s$.
- не учитываются значения расстояний

Пример зависимости LOO от числа соседей

Пример. Задача Iris, усреднение по 50 случайным разбиениям

частота ошибок



— несмещённое число ошибок LOO

— смещённое число ошибок, когда объект учитывается как сосед самого себя

Метод окна Парзена

$w(i, x) = K\left(\frac{\rho(x, x^{(i)})}{h}\right)$, где h — ширина окна,
 $K(r)$ — ядро, не возрастает и положительно на $[0, 1]$.

Метод парзеновского окна *заданной ширины* h :

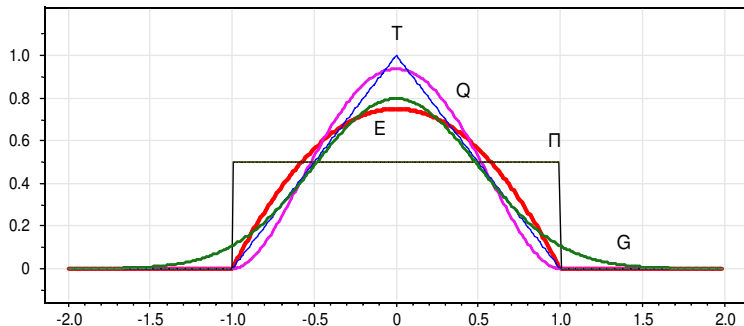
$$a(x; X^\ell, h, K) = \arg \max_{y \in Y} \sum_{i=1}^{\ell} [y_i = y] K\left(\frac{\rho(x, x_i)}{h}\right)$$

Метод парзеновского окна *переменной ширины* $h = \rho(x, x^{(k+1)})$:

$$a(x; X^\ell, k, K) = \arg \max_{y \in Y} \sum_{i=1}^{\ell} [y_i = y] K\left(\frac{\rho(x, x_i)}{\rho(x, x^{(k+1)})}\right)$$

Оптимизация параметров — по критерию LOO:

- выбор ширины окна h или числа соседей k
- выбор ядра K слабо влияет на качество классификации

Часто используемые ядра $K(r)$ 

$P(r) = [|r| \leq 1]$ — прямоугольное

$T(r) = (1 - |r|)[|r| \leq 1]$ — треугольное

$E(r) = (1 - r^2)[|r| \leq 1]$ — квадратичное (Епанечникова)

$Q(r) = (1 - r^2)^2[|r| \leq 1]$ — четвертое

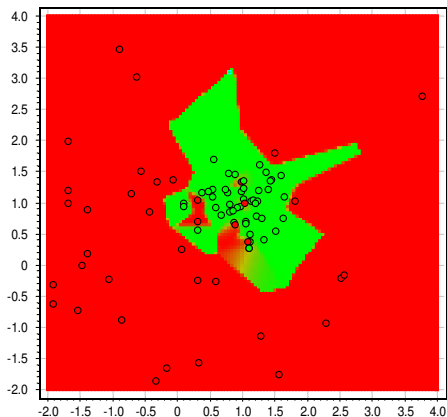
$G(r) = \exp(-2r^2)$ — гауссовское

Парзеновское окно ширины h

Пример: двумерная выборка, два класса $Y = \{-1, +1\}$.

$$a(x) = \arg \max_{y \in Y} \Gamma_y(x) = \text{sign}(\underbrace{\Gamma_{+1}(x) - \Gamma_{-1}(x)}})$$

$h = 0.05$

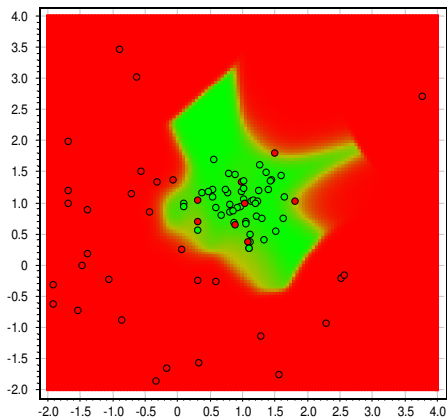


Парzenовское окно ширины h

Пример: двумерная выборка, два класса $Y = \{-1, +1\}$.

$$a(x) = \arg \max_{y \in Y} \Gamma_y(x) = \text{sign}(\underbrace{\Gamma_{+1}(x) - \Gamma_{-1}(x)}})$$

$h = 0.2$

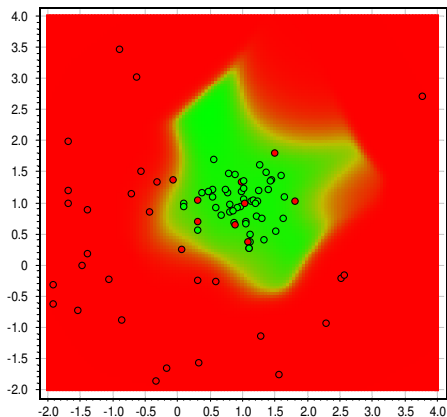


Парzenовское окно ширины h

Пример: двумерная выборка, два класса $Y = \{-1, +1\}$.

$$a(x) = \arg \max_{y \in Y} \Gamma_y(x) = \text{sign}(\underbrace{\Gamma_{+1}(x) - \Gamma_{-1}(x)}})$$

$h = 0.3$

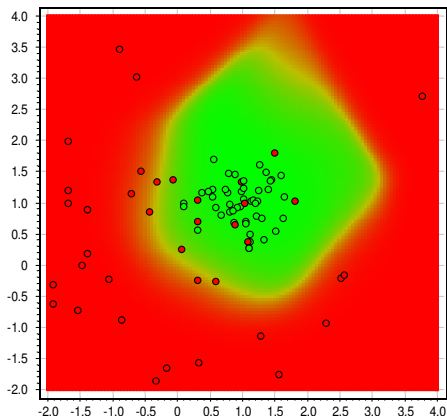


Парzenовское окно ширины h

Пример: двумерная выборка, два класса $Y = \{-1, +1\}$.

$$a(x) = \arg \max_{y \in Y} \Gamma_y(x) = \text{sign}(\underbrace{\Gamma_{+1}(x) - \Gamma_{-1}(x)}})$$

$h = 0.5$

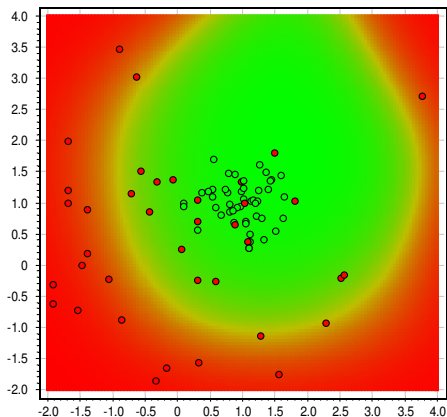


Парzenовское окно ширины h

Пример: двумерная выборка, два класса $Y = \{-1, +1\}$.

$$a(x) = \arg \max_{y \in Y} \Gamma_y(x) = \text{sign}(\underbrace{\Gamma_{+1}(x) - \Gamma_{-1}(x)}})$$

$h = 1.0$

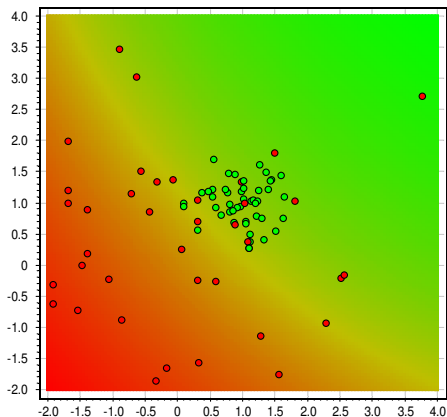


Парzenовское окно ширины h

Пример: двумерная выборка, два класса $Y = \{-1, +1\}$.

$$a(x) = \arg \max_{y \in Y} \Gamma_y(x) = \text{sign}(\underbrace{\Gamma_{+1}(x) - \Gamma_{-1}(x)}_{\text{разность}})$$

$h = 5.0$



Метод потенциальных функций

$$w(i, x) = \gamma^{(i)} K\left(\frac{\rho(x, x^{(i)})}{h^{(i)}}\right)$$

Более простая запись (без ранжирования объектов):

$$a(x; X^\ell) = \arg \max_{y \in Y} \sum_{i=1}^{\ell} [y_i = y] \gamma_i K\left(\frac{\rho(x, x_i)}{h_i}\right),$$

где γ_i — веса объектов, $\gamma_i \geq 0$, $h_i > 0$.

Физическая аналогия:

γ_i — величина «заряда» в точке x_i ;

h_i — «радиус действия» потенциала с центром в точке x_i ;

y_i — знак «заряда» (в случае двух классов $Y = \{-1, +1\}$);

в электростатике $K(r) = \frac{1}{r}$ или $\frac{1}{r+a}$,

для задач классификации нет таких ограничений на K .

Метод потенциальных функций = линейный классификатор

Два класса: $Y = \{-1, +1\}$.

$$\begin{aligned} a(x; X^\ell) &= \arg \max_{y \in Y} \Gamma_y(x) = \text{sign}(\Gamma_{+1}(x) - \Gamma_{-1}(x)) = \\ &= \text{sign} \sum_{i=1}^{\ell} \gamma_i y_i K\left(\frac{\rho(x, x_i)}{h_i}\right). \end{aligned}$$

Сравним с линейной моделью классификации:

$$a(x) = \text{sign} \sum_{j=1}^n \gamma_j f_j(x).$$

- функции $f_j(x) = y_j K\left(\frac{1}{h_j} \rho(x, x_j)\right)$ — признаки объекта x
- γ_j — веса линейного классификатора
- $n = \ell$ — число признаков равно числу объектов обучения

Задача регрессии и метод наименьших квадратов

- $X = \mathbb{R}^n$ — объекты; $Y = \mathbb{R}$ — ответы;
 $X^\ell = (x_i, y_i)_{i=1}^\ell$ — обучающая выборка;
 $y_i = f(x_i)$, $f: X \rightarrow Y$ — неизвестная зависимость;
- $a(x) = g(x, \alpha)$ — параметрическая модель зависимости,
 $\alpha \in \mathbb{R}^p$ — вектор параметров модели.

- Метод наименьших квадратов (МНК):

$$Q(\alpha, X^\ell) = \sum_{i=1}^{\ell} w_i (g(x_i, \alpha) - y_i)^2 \rightarrow \min_{\alpha},$$

где w_i — вес, степень важности i -го объекта.

- **Недостаток:**
надо иметь хорошую параметрическую модель $g(x, \alpha)$

Непараметрическая регрессия. Формула Надарая–Ватсона

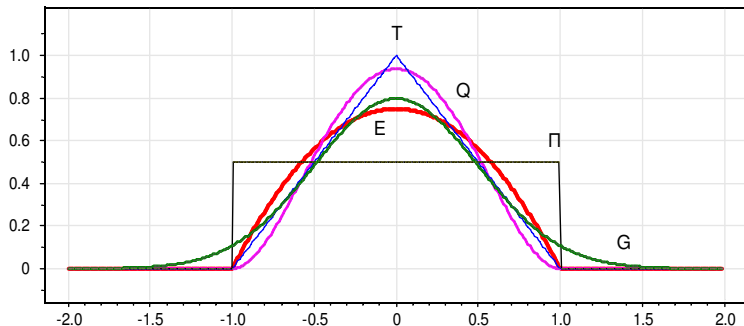
Приближение константой $g(x, \alpha) = \alpha$ в окрестности $x \in X$:

$$Q(\alpha; X^\ell) = \sum_{i=1}^{\ell} w_i(x) (\alpha - y_i)^2 \rightarrow \min_{\alpha \in \mathbb{R}}$$

где $w_i(x) = K\left(\frac{\rho(x, x_i)}{h}\right)$ — веса объектов x_i относительно x ;
 $K(r)$ — ядро, невозрастающее, ограниченное, гладкое;
 h — ширина окна сглаживания.

Формула ядерного сглаживания Надарая–Ватсона:

$$a_h(x; X^\ell) = \frac{\sum_{i=1}^{\ell} y_i w_i(x)}{\sum_{i=1}^{\ell} w_i(x)} = \frac{\sum_{i=1}^{\ell} y_i K\left(\frac{\rho(x, x_i)}{h}\right)}{\sum_{i=1}^{\ell} K\left(\frac{\rho(x, x_i)}{h}\right)}.$$

Часто используемые ядра $K(r)$ — те же, что в методе Парзена

$P(r) = [|r| \leq 1]$ — прямоугольное

$T(r) = (1 - |r|) [|r| \leq 1]$ — треугольное

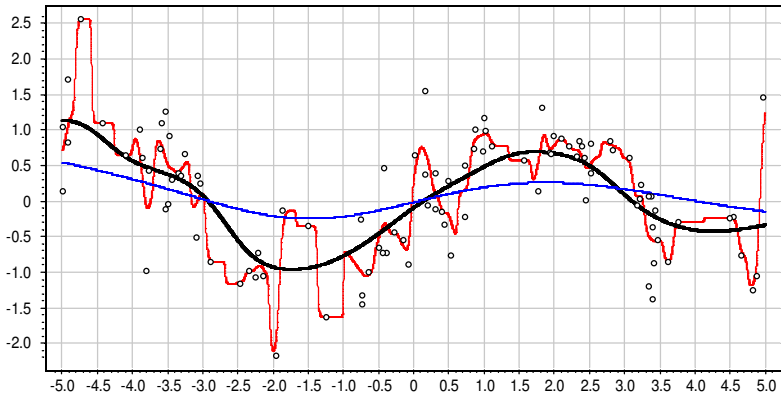
$E(r) = (1 - r^2) [|r| \leq 1]$ — квадратичное (Епанечникова)

$Q(r) = (1 - r^2)^2 [|r| \leq 1]$ — четвертое

$G(r) = \exp(-2r^2)$ — гауссовское

Выбор ядра K и ширины окна h

$h \in \{0.1, 1.0, 3.0\}$, гауссовское ядро $K(r) = \exp(-2r^2)$

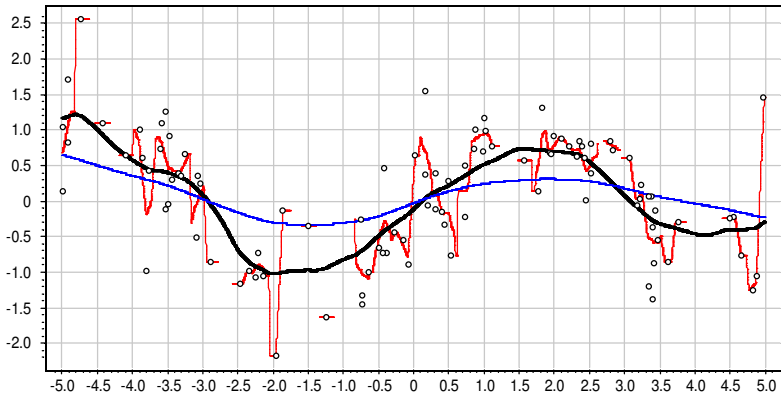


Гауссовское ядро \Rightarrow гладкая аппроксимация

Ширина окна существенно влияет на точность аппроксимации

Выбор ядра K и ширины окна h

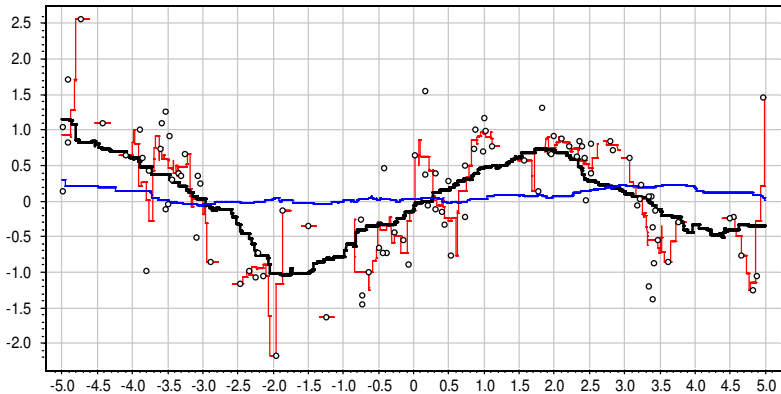
$h \in \{0.1, 1.0, 3.0\}$, треугольное ядро $K(r) = (1 - |r|) [|r| \leq 1]$



Треугольное ядро \Rightarrow кусочно-линейная аппроксимация
Аппроксимация не определена, если в окне нет точек выборки

Выбор ядра K и ширины окна h

$h \in \{0.1, 1.0, 3.0\}$, прямоугольное ядро $K(r) = [|r| \leq 1]$



Прямоугольное ядро \Rightarrow кусочно-постоянная аппроксимация
Выбор ядра слабо влияет на точность аппроксимации

Выбор ядра K и ширины окна h

- Ядро $K(r)$
 - существенно влияет на гладкость функции $a_h(x)$,
 - слабо влияет на качество аппроксимации.
- Ширина окна h
 - существенно влияет на качество аппроксимации.
- Переменная ширина окна $h(x)$ по k ближайшим соседям:

$$w_i(x) = K\left(\frac{\rho(x, x_i)}{h(x)}\right),$$

где $h(x) = \rho(x, x^{(k+1)})$, $x^{(k)}$ — k -й сосед объекта x .

- Оптимизация h (или k) по скользящему контролю:

$$\text{LOO}(h, X^\ell) = \sum_{i=1}^{\ell} \left(a_h(x_i; X^\ell \setminus \{x_i\}) - y_i \right)^2 \rightarrow \min_h.$$

Постановка задачи кластеризации

Дано:

$X^\ell = \{x_1, \dots, x_\ell\}$ — обучающая выборка;

$\rho(x_i, x_s)$ — функция расстояния между объектами.

Найти:

$y_i \in Y$ — метки кластеров объектов:

— каждый кластер состоит из близких объектов;

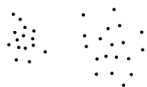
— объекты разных кластеров существенно различны.

Кластеризация — это *обучение без учителя*.

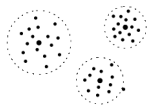
Решение задачи кластеризации принципиально неоднозначно:

- различные критерии качества кластеризации
- различные эвристические методы кластеризации
- различные варианты функции расстояния ρ

Типы кластерных структур



внутрикластерные расстояния, как правило, меньше межкластерных



кластеры с центром



кластеры могут соединяться перемычками



кластеры могут накладываться на разреженный фон из редко расположенных объектов

Типы кластерных структур



ленточные кластеры



перекрывающиеся кластеры



кластеры могут образовываться не по сходству, а по иным типам регулярностей



кластеры могут вообще отсутствовать

Метод k -средних (k -means)

Объекты x_i задаются векторами признаков (x_i^1, \dots, x_i^n) .

Вход: X^ℓ — обучающая выборка, параметр k ;

Выход: центры кластеров μ_y , $y \in Y$ и кластеризация $y_i \in Y$;

1 начальное приближение центров μ_y , $y \in Y$;

2 **повторять**

3 | отнести каждый x_i к ближайшему центру:

$$y_i := \arg \min_{y \in Y} \rho(x_i, \mu_y), \quad i = 1, \dots, \ell;$$

4 | вычислить новые положения центров:

$$\mu_{yj} := \frac{\sum_{i=1}^{\ell} [y_i = y] x_i^j}{\sum_{i=1}^{\ell} [y_i = y]}, \quad y \in Y, \quad j = 1, \dots, n;$$

5 **пока** y_i не перестанут изменяться;

Пример. Результат кластеризации методом k -средних

K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross



Агломеративная иерархическая кластеризация

Алгоритм Ланса-Уильямса [1967] основан на оценивании расстояний $R(U, V)$ между парами кластеров U, V .

- 1 сначала все кластеры одноэлементные: $C_1 := \{\{x_1\}, \dots, \{x_\ell\}\}$;
- 2 расстояния между ними: $R(\{x_i\}, \{x_s\}) := \rho(x_i, x_s)$;
- 3 **для всех** $t = 2, \dots, \ell$ (t — номер итерации):
 - 4 найти в C_{t-1} два ближайших кластера:
 $(U, V) := \arg \min_{U \neq V} R(U, V)$;
 - 5 $R_t := R(U, V)$;
 - 6 слить их в один кластер $W := U \cup V$;
 $C_t := C_{t-1} \cup \{W\} \setminus \{U, V\}$;
 - 7 **для всех** $S \in C_t$
└ вычислить $R(W, S)$ по формуле Ланса-Уильямса;

Формула Ланса-Уильямса

Как определить расстояние $R(W, S)$
между кластерами $W = U \cup V$ и S ,
зная расстояния $R(U, S)$, $R(V, S)$, $R(U, V)$?

Формула, обобщающая большинство разумных способов
определить это расстояние [Ланс, Уильямс, 1967]:

$$\begin{aligned} R(U \cup V, S) = & \alpha_U \cdot R(U, S) + \\ & + \alpha_V \cdot R(V, S) + \\ & + \beta \cdot R(U, V) + \\ & + \gamma \cdot |R(U, S) - R(V, S)|, \end{aligned}$$

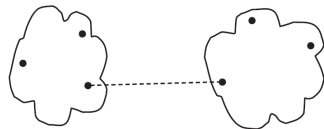
где α_U , α_V , β , γ — числовые параметры.

Частные случаи формулы Ланса-Уильямса

1. Расстояние ближнего соседа:

$$R^b(W, S) = \min_{w \in W, s \in S} \rho(w, s);$$

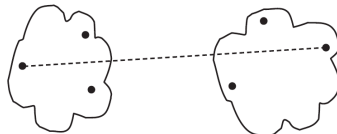
$$\alpha_U = \alpha_V = \frac{1}{2}, \quad \beta = 0, \quad \gamma = -\frac{1}{2}.$$



2. Расстояние дальнего соседа:

$$R^a(W, S) = \max_{w \in W, s \in S} \rho(w, s);$$

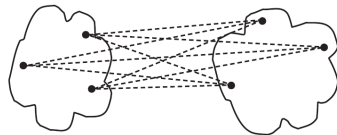
$$\alpha_U = \alpha_V = \frac{1}{2}, \quad \beta = 0, \quad \gamma = \frac{1}{2}.$$



3. Групповое среднее расстояние:

$$R^r(W, S) = \frac{1}{|W||S|} \sum_{w \in W} \sum_{s \in S} \rho(w, s);$$

$$\alpha_U = \frac{|U|}{|W|}, \quad \alpha_V = \frac{|V|}{|W|}, \quad \beta = \gamma = 0.$$



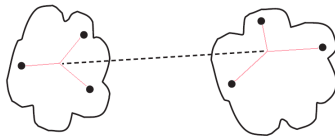
Частные случаи формулы Ланса-Уильямса

4. Расстояние между центрами:

$$R^4(W, S) = \rho^2 \left(\sum_{w \in W} \frac{w}{|W|}, \sum_{s \in S} \frac{s}{|S|} \right);$$

$$\alpha_U = \frac{|U|}{|W|}, \quad \alpha_V = \frac{|V|}{|W|},$$

$$\beta = -\alpha_U \alpha_V, \quad \gamma = 0.$$



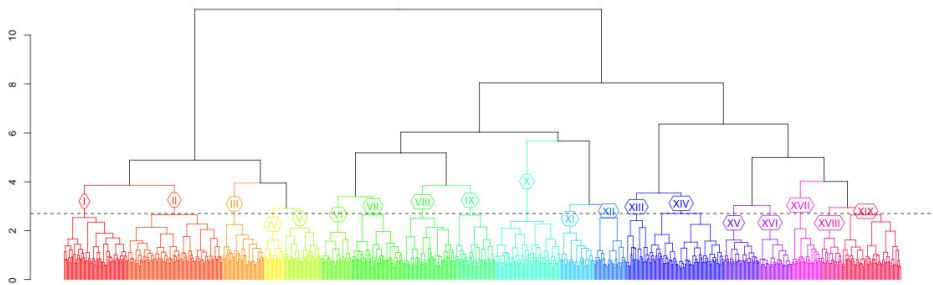
5. Расстояние Уорда — рекомендуется для использования

$$R^y(W, S) = \frac{|S||W|}{|S|+|W|} \rho^2 \left(\sum_{w \in W} \frac{w}{|W|}, \sum_{s \in S} \frac{s}{|S|} \right);$$

$$\alpha_U = \frac{|S|+|U|}{|S|+|W|}, \quad \alpha_V = \frac{|S|+|V|}{|S|+|W|}, \quad \beta = \frac{-|S|}{|S|+|W|}, \quad \gamma = 0.$$

Дендрограмма — визуализация иерархической кластеризации

- По вертикальной оси откладываются расстояния R_t
- Уровень отсечения определяет число кластеров
- Дендрограмма не имеет самопересечений и группирует объекты в кластеры вдоль горизонтальной оси



Ещё один способ визуализации иерархической кластеризации



Резюме

- Функции расстояния полезны для решения задач классификации, регрессии, кластеризации.
- Ширину окна h оптимизируют по критерию LOO:
 - при $h \rightarrow 0$ модель усложняется, начинает настраиваться на шум и переобучается;
 - при $h \rightarrow \infty$ модель упрощается, начинает стремиться к константе и недообучается;
 - «золотую середину» находят перебором h по сетке.
- Кластеризация неоднозначна, результат может зависеть от метрики, от критерия, от метода
- Кластеризация используется для выделения структурных особенностей выборки объектов

Воронцов Константин Вячеславович

voron@forecsys.ru

www.MachineLearning.ru • Участник:Vokov

Если что-то было не понятно,
не стесняйтесь подходить и спрашивать :)