

Метод k ближайших соседей

М. Корневская, гр. 23641/3

18 октября 2017 г.

- Обучение с учителем
 - Классификация
 - Регрессия
 - Ранжирование
 - Прогнозирование
- Обучение без учителя
- Частичное обучение
и др.

Метод k ближайших соседей использует расстояния (метрики) в пространстве объектов.

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 $\rho : X \times X \rightarrow [0, \infty)$.

Пример. Евклидово расстояние и его обобщение:

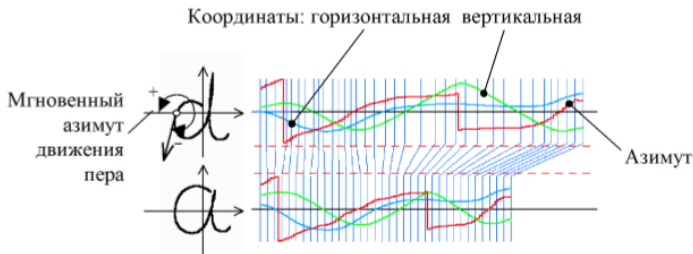
$$\rho(x, x_i) = \left(\sum_{j=1}^n |x^j - x_i^j|^2 \right)^{1/2} \quad \rho(x, x_i) = \left(\sum_{j=1}^n w_j |x^j - x_i^j|^p \right)^{1/p}$$

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

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

Другие примеры расстояний

- между текстами (редакционное расстояние Левенштейна):
GCTAAAGGTACGCC . . TTTAGAAA . GGGCCATTAGGAAA TTGC
GACTAA AGCCTATTTACAAATGGGCCATTAGG . . . TTGC
- между сигналами (энергия сжатий и растяжений):



Обобщенный метрический классификатор

Для произвольного $x \in X$ отранжируем объекты x_1, \dots, x_l :

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

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

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

Метрический алгоритм классификации:

$$a(x, X^l) = \arg \max_{y \in Y} \sum_{i=1}^l [y^{(i)} = y] w(i, x),$$

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

$\Gamma_y(x) = \sum_{i=1}^l [y^{(i)} = y] w(i, x)$ — оценка близости объекта x к классу y .

Метод k ближайших соседей (k nearest neighbors, kNN)

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

Преимущества:

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

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

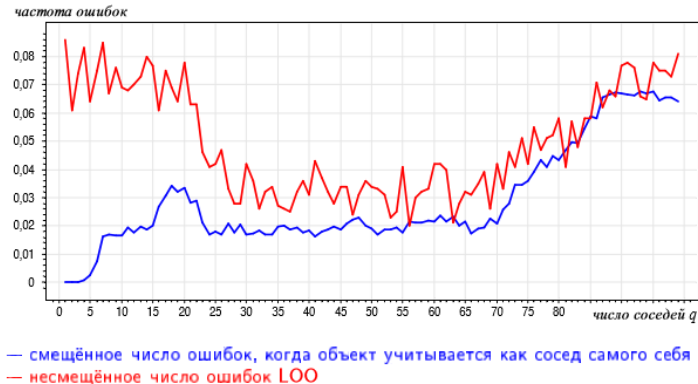
Проблемы:

- возможны ситуации, когда классификация не однозначна:
 $\Gamma_y(x) = \Gamma_s(x)$ для пары классов $y \neq s$
- учитываются не значения расстояний, а только их ранги

Пример зависимости LOO от числа соседей

Пример.

Задача Iris.



В реальных задачах минимум редко бывает при $k = 1$.