

Лекция kNN деревья решения
 1. Точные методы

kNN - сложность $O(ld)$

Если $\vec{x} \in \mathbb{R}^d$, $m \in \mathbb{N}$, то $\vec{x} \Rightarrow \vec{z} \in \mathbb{R}^m$

$$d(\vec{u}, \vec{x}, \rho) \Rightarrow a(\vec{u}, \vec{z}, \rho) \quad O(lm)$$

k_d - деревья

$$d \approx 10 \sim 20 \quad O(\log l)$$

$$d > 20 \quad O(ld)$$

2. Приближенные методы

1. Заполняем не всю обучающую выборку
 $|\{x_i\}| = l$, а лишь часть

STOLP

2. Усечь k ближайших соседей к ближайшему

Опр 2.1

1. Если $\rho(x, y) \leq d_1 \Rightarrow P_{f \in F} [f(x) = f(y)] \geq p_1$

2. Если $\rho(x, y) \not\leq d_2 \Rightarrow P_{f \in F} [f(x) = f(y)] \leq p_2$

где $d_1 \leq d_2$ $p_1 \geq p_2$

Пример

$$p = \frac{|S(A \cap B)|}{|S(A \cup B)|}$$

$U = \{u_1, \dots, u_n\}$ - универсум

π - перестановка $\pi \in \text{IDP?}$

$$f_\pi(A) = \min \{ \pi(i) \mid u_i \in A \}$$

U - это слова

π - степень важности слов, чем меньше

$\pi(i)$, тем важнее i -е слово

Примерный малый уровень важности

Для неидеальных текстов, слово

сифрофразатрон

Покажем, что множество всех
MinKash-функций $F = \{f_\pi \mid \pi \in \text{Sym}(U)\}$
является (d_1, d_2, p_1, p_2) .

Сначала докажем:

Вероятность того что случайно выбранная
функция $f_\pi \in F$ будет принимать
одинаковые значения на двух заданных
множествах A и B , равна котф-гу

Отсюда
$$\frac{|A \cap B|}{|A \cup B|}$$

$$1) u \in A, u \in B - p$$

$$2) u \in A, u \notin B \text{ или } u \in B, u \notin A - q$$

$$3) u \notin A, u \notin B - o$$

$$\frac{p}{p+q} - \text{вероятность}$$

$$\text{Лемма } p_J(A, B) \leq d_1 \Rightarrow 1 - p_J(A, B) \geq 1 - d_1$$

$$\Rightarrow p_1 \geq 1 - d_1 \Rightarrow$$

$$(d_1, d_2, 1 - d_1, 1 - d_2)$$

Композиция хэш функций

f - хэш функция

T - хэш таблица

разместим \vec{x}_i в ячейку с номером $f(x)$

При этом требуется пойти к элементам n и r

Однако p_1 и p_2 мало отличаются
потому что k -потных соседей
либо много либо мало

$$g_1(x) = (f_{11}(x), \dots, f_{1m}(x)) - T_1$$

$$g_L(x) = (f_{L1}(x), \dots, f_{Lm}(x)) - T_L$$

Чтобы найти k ближайших соседей
 для нового объекта x , выберем объекты
 из тех $g_1(x), \dots, g_L(x)$ таблиц T_1, \dots, T_L
 и вернем k наиболее близких x из них

Два базовых параметра

число m — кол-во функций в одной таблице

число L — кол-во таблиц

$$m \nearrow \Rightarrow P(f(\vec{x}) = f(\vec{y})) \nearrow \text{ где } \vec{x} \neq \vec{y}$$

но если $m \gg m_0$, то схожих объектов
 вообще не будет

$$L \nearrow \Rightarrow P(f(\vec{x}) = f(\vec{y})) \nearrow \text{ где } \rho(\vec{x}, \vec{y}) < d_1$$

но если $L \gg L_0$, то будет слишком
 много индифферентов

$$(d_1, d_2, 1 - (1 - p_1^m)^L, 1 - (1 - p_2^m)^L)$$

Если расстояние между этими объектами
 велико т.е. $\rho(\vec{x}, \vec{y}) \geq d_2 + \epsilon$ вероятности
 совпадения m базовых функций p_2^m

Теоретические гарантии

Будем говорить, что алгоритм решает задачу поиска ϵ -ближайшего соседа

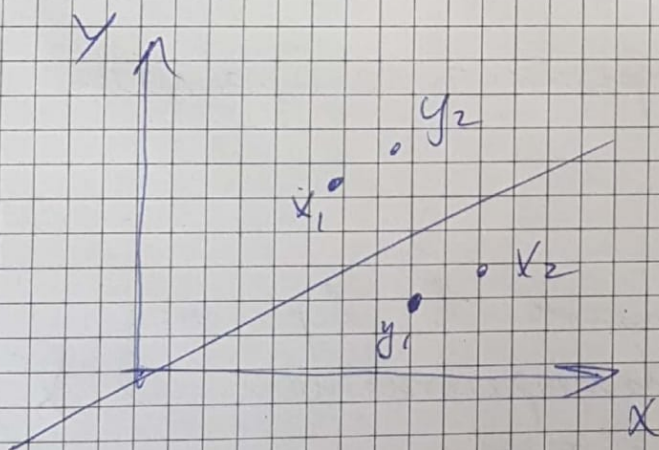
если для u он с $p = 1 - \epsilon$ возвращает объект выбора удаленный от u не более чем в ϵ раз

$O(d^r \log L)$, где r для многих функций разложения имеет порядок $\frac{1}{\epsilon}$

Хэм-функция для косинусного разложения

$$H = \{f_w(\vec{x}) = \text{sign}(\langle \vec{w}, \vec{x} \rangle) \mid w \in \mathbb{R}^d\}$$

Покажите, что Хэм-ф. косинуса является (d_1, d_2, p_1, p_2) - чувствительными



$$\angle \vec{x}, \vec{y} = \theta$$

$$\frac{180 - \theta}{180}$$

$$(d_1, d_2, 1 - \frac{d_1}{\pi}, 1 - \frac{d_2}{\pi})$$

Хэм евклидовой метрики

$$f_{w,b}(x) = \left\lfloor \frac{\langle w, x \rangle + b}{r} \right\rfloor \quad (w \in \mathbb{R}^d, b \in [0, r])$$

Дерево

$$\forall \text{ node} : \beta_v : X \mapsto \{0, 1\}$$

$$\forall \text{ term} : c_v \in Y$$

Пример с кредитом

1) одномерные предикаты

2) линейные предикаты

$$\beta_v(x) = [\langle \vec{w}, \vec{x} \rangle < S]$$

$$3) \beta_v(x) = [p(x, x_v) < S] \quad x_v - \text{этalon}$$

одномерные предикаты с произвольной размерностью

Алгоритм:

- выдел предиката в вершине
- критерии информативности $\Phi(X, \beta_v)$
- критерии останова
- метод обр-ки пр-ых значений
- метод стрижки

R_m — мн-во объектов об-го выбора
попадающих в вершину m

$$N_m = |R_m|$$

R_{mk} — доля объектов класса k
 $k \in \{1, \dots, K\}$ в вершине m .

$$R_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} [y_i = k]$$

K_m — обозначим класса ~~сбих~~
присущих объектам ~~большин~~

$$K_m = \arg \max_k R_{mk}$$