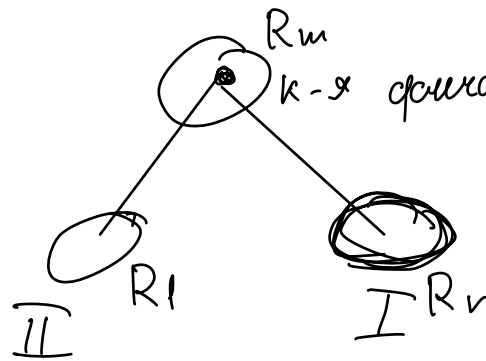


Выбор предикатов

X



k-я функция, τ

$x^k \geq \tau$ — левое

$x^k < \tau$ — правое

R_m — мн-во объектов

$$Q(R_m) = H(R_m) - \frac{|R_l|}{|R_m|} H(R_l) - \frac{|R_r|}{|R_m|} H(R_r)$$

$H(R)$ — критерий информативности.

$$H(R) \rightarrow \min$$

$$Q(R_m) \rightarrow \max$$

$$H(R) = \min_{c \in Y} \frac{1}{|R|} \sum_{(x,y) \in R} L(y, c)$$

Регрессия

$$H(R) = \min_{c \in Y} \frac{1}{|R|} \sum_{(x,y) \in R} (y - c)^2$$

$$H(R) = \min_{c \in Y} \frac{1}{|R|} \sum_{(x,y) \in R} (y - \bar{y})^2, \quad \bar{y} = \frac{1}{|R|} \sum y$$

Классификация

$$R, \quad k \in \{1 \dots K\}$$

P_k - доля объектов класса k в R

$$P_k = \frac{1}{|R|} \sum_{(x,y) \in R} [y_i = k]$$

k^* - наиболее мощный класс

$$k^* = \arg \max_k P_k$$

$$H(R) = \min_{c \in Y} \frac{1}{|R|} \sum_{(x,y) \in R} [y_i \neq c] \quad \textcircled{=}$$

$$c = P_{k^*}$$

$$\textcircled{=} 1 - P_{k^*}$$

Критерий Дженсена

$$c = (c_1, \dots, c_K), \quad \sum_{k=1}^K c_k = 1$$

Критерий Брера

$$H(R) = \min_{\sum_k c_k = 1} \frac{1}{|R|} \sum_{(x,y) \in R} \sum_{k=1}^K (c_k - \underbrace{[y_i = k]}_?)^2$$

$$C^* = (p_1, \dots, p_K)$$

$$H(R) = \sum_{k=1}^K p_k (1 - p_k) = \sum_{k=1}^K (p_k - p_k^2) =$$

$$= \underbrace{\sum_{k=1}^K p_k}_1 - \sum_{k=1}^K p_k^2$$

Задача 1. $H(R) = \sum_{k \neq k'} p_k p_{k'} = \sum_{k=1}^K \sum_{k' \neq k}^K p_k p_{k'} =$

$$= \sum_{k=1}^K p_k \left(\sum_{k' \neq k}^K p_{k'} \right) = \sum_{k=1}^K p_k (1 - p_k)$$

Задача 2: Есть вершины m . Объекты R
 $a(x)$, которая выбирает класс случайно.

Но класс k выбирается с вер-ю p_k

E (частота ошибок) равно индексу Джинни.

$$E = \frac{1}{|R|} \sum_{x,y} [y \neq a(x)] = \frac{1}{|R|} \sum_{x,y} E[y \neq a(x)] =$$

$$= \frac{1}{|R|} \sum_{x,y} (1 - p_y) \odot$$

$$\begin{aligned} \textcircled{=}& \frac{1}{|R|} \sum_{x,y} \sum_{k=1}^K [y_i = k] (1 - p_k) = \\ &= \sum_{k=1}^K \underbrace{\frac{\sum_{x,y} [y_i = k]}{|R|}}_{p_k} (1 - p_k) = \sum_{k=1}^K p_k (1 - p_k) \end{aligned}$$

$$Q(R_m) = \cancel{H(R_m)} - \underbrace{\frac{|R_l|}{|R_m|}}_{\text{}} H(R_l) - \underbrace{\frac{|R_r|}{|R_m|}}_{\text{}} H(R_r) \rightarrow \max \textcircled{=}$$

$$\textcircled{=} -\frac{1}{|R_m|} \left(|R_l| H(R_l) + |R_r| H(R_r) \right) \textcircled{=} \text{ } \begin{matrix} p_{mk} - \text{вероятность} \\ \text{в вершине } m \end{matrix}$$

$$\begin{aligned} \textcircled{=}& -\frac{1}{|R_m|} \left(\cancel{|R_l|} - |R_l| \left(\sum_{k=1}^K p_{lk}^2 \right) + \cancel{|R_r|} - |R_r| \left(\sum_{k=1}^K p_{rk}^2 \right) \right) = \\ & |R_l| \left(\sum_{k=1}^K p_{lk}^2 \right) + |R_r| \left(\sum_{k=1}^K p_{rk}^2 \right) \end{aligned}$$

число пар (x_i, x_j) т.ч. оба объекта
находятся в одних и тех же подгруппах и
 $y_i = y_j$

$$p_{lk} |R_l| \left(\sum_{k=1}^K p_{lk}^2 \right) |R_l| \textcircled{2} + \left(\sum_{k=1}^K p_{rk}^2 \right) |R_r| \textcircled{2}$$

Эмпирический критерий

$$\prod_{k=1}^K P_k^{[y_i=k]}$$

$$H(R) = \min_{\sum c_k = 1} \left(- \frac{1}{|R|} \sum_{x,y} \sum_{k=1}^K [y_i=k] \log c_k \right) \rightarrow \min_{c_k}$$

$$L(c, \lambda) = - \frac{1}{|R|} \sum_{x,y} \sum_{k=1}^K [y_i=k] \log c_k + \lambda \sum_{k=1}^K c_k \rightarrow \min_{c_k}$$

$$\frac{\partial L}{\partial c_k} = - \frac{1}{|R|} \sum_{x,y} [y_i=k] \frac{1}{c_k} + \lambda \quad \frac{P_k}{c_k} + \lambda = 0$$

$$c_k = P_k / \lambda$$

$$1 = \frac{1}{\lambda} \sum_{k=1}^K P_k = \frac{1}{\lambda} \Rightarrow \lambda = 1 \quad c_k = P_k$$

$$H(R) = - \sum_{k=1}^K P_k \log P_k$$

$\log_2 K$