

4 Теоретические задачи

4.1 наивный байес и центроидный классификатор

$$\begin{aligned} \underset{y}{\operatorname{argmax}} \left(\prod_{k=1}^n \frac{1}{2\pi\sigma^2} e^{-\frac{(x^{(k)} - \mu_{y_k})^2}{2\sigma^2}} Pr(y) \right) = \\ = \underset{y}{\operatorname{argmin}} \left(\sum_{k=1}^n (x^{(k)} - \mu_{y_k})^2 \right) \end{aligned}$$

Поскольку $Pr(y) = \text{const}$, а максимум оставшейся функции является минимумом аргумента экспоненты, в которой также можно избавиться от константного множителя

4.2 ROC-AUC случайных ответов Пусть количество 1 в выборке равно k , размер выборки n , тогда:

$$\begin{aligned} \xi_i &= \text{Bern}(p) \\ tp &= \sum_{i=1}^k \xi_i, fp = \sum_{i=1}^{n-k} \xi_i \\ tn &= \sum_{i=1}^{n-k} 1 - \xi_i, fn = \sum_{i=1}^k 1 - \xi_i \end{aligned}$$

Получаем искомые случайные величины:

$$\begin{aligned} FPR &= \frac{fp}{fp + tn} = \frac{\sum_{i=1}^{n-k} \xi_i}{\sum_{i=1}^{n-k} \xi_i + \sum_{i=1}^{n-k} (1 - \xi_i)} = \frac{\sum_{i=1}^{n-k} \xi_i}{n - k} \\ TPR &= \frac{tp}{tp + fn} = \frac{\sum_{i=1}^k \xi_i}{\sum_{i=1}^k \xi_i + \sum_{i=1}^k (1 - \xi_i)} = \frac{\sum_{i=1}^k \xi_i}{k} \end{aligned}$$

взяв математическое ожидание, получим:

$$\mathbb{E}FPR = p$$

$$\mathbb{E}TPR = p$$

Значит средняя точка лежит на диагонали квадрата. Значит площадь в среднем равна $\frac{1}{2}$

4.3 Ошибка 1NN оптимального байесовского классификатора

$$\begin{aligned} \lim_{n \rightarrow \infty} Pr(y \neq y_n) &= \lim_{n \rightarrow \infty} \sum_{y \neq y_n \in \{0,1\}} Pr(y, y_n | x, x_n) = \\ &= \lim_{n \rightarrow \infty} \sum_{y \neq y_n \in \{0,1\}} Pr(y|x) Pr(y_n|x_n) = \sum_{y \in \{0,1\}} Pr(y|x) (1 - Pr(y|x)) = \\ &= Pr(0|x)(1 - Pr(0|x)) + Pr(1|x)(1 - Pr(1|x)) = \\ &= 2 \max_{y \in \{0,1\}} Pr(y|x) (1 - \max_{y \in \{0,1\}} Pr(y|x)) \leq 2E_B = \\ &= 2(1 - \max_{y \in \{0,1\}} Pr(y|x)) \end{aligned}$$