

# Теоретические задачи

## 1.1 Ответы в листьях регрессионного дерева

Рассмотрим некоторый лист дерева. Пусть в него попали объекты  $y_1, \dots, y_n$ . Посчитаем матожидание MSE для некоторого объекта  $y$  в первом и во втором случае:

$$MSE = E(y - \hat{y})^2 = Ey^2 + E\hat{y}^2 - 2Ey \cdot \hat{y} = Ey^2 + E\hat{y}^2 - 2Ey \cdot E\hat{y}$$

Значит, разность ошибок:

$$err_1 - err_2 = E\hat{y}_1^2 - E\hat{y}_2^2 - 2Ey \cdot (E\hat{y}_1 - E\hat{y}_2)$$

$$1) \hat{y}_1 = \frac{\sum_{i=1}^n y_i}{n}$$

$$E\hat{y}_1 = \frac{\sum_{i=1}^n Ey_i}{n} = \frac{\sum_{i=1}^n y_i}{n}$$

$$E\hat{y}_1^2 = E\left(\frac{\sum_{i=1}^n y_i}{n}\right)^2$$

$$2) E\hat{y}_2 = \frac{\sum_{i=1}^n y_i}{n}$$

$$E\hat{y}_2^2 = \frac{\sum_{i=1}^n y_i^2}{n}$$

Разность ошибок:

$$err_1 - err_2 = E\left(\frac{\sum_{i=1}^n y_i}{n}\right)^2 - \frac{\sum_{i=1}^n y_i^2}{n} - 2Ey \cdot \left(\frac{\sum_{i=1}^n y_i}{n} - \frac{\sum_{i=1}^n y_i}{n}\right) \leq 0$$

Использовано неравенство Коши-Буняковского  $\left(\sum_{i=1}^n y_i\right)^2 \leq \left(\sum_{i=1}^n y_i^2\right) \cdot \left(\sum_{i=1}^n 1^2\right) = n \cdot \sum_{i=1}^n y_i^2$

## 1.3 Unsupervised decision tree

Энтропия непрерывного распределения выражается через его плотность по формуле:  $-\int \dots \int f(x) \cdot \ln(f(x))$ , где  $f(x)$  - плотность распределения. Плотность многомерного нормального распределения:  $f(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(x-a)^T \Sigma^{-1}(x-a)}$ .

Заметим, что  $\frac{1}{2}(x-a)^T \Sigma^{-1}(x-a) = \sum_i \sum_j (x_i - a_i) \cdot \Sigma_{i,j}^{-1} \cdot (x_j - a_j)$  (в силу произведения матриц).

Найдем энтропию нормального распределения:

$$\begin{aligned} H &= -\int \dots \int f(x) \cdot \left(-\frac{1}{2}(x-a)^T \Sigma^{-1}(x-a) - \frac{1}{2} \cdot \ln((2\pi)^n |\Sigma|)\right) = \frac{1}{2} E\left(\sum_i \sum_j (x_i - a_i) \cdot \Sigma_{i,j}^{-1} \cdot (x_j - a_j)\right) + \\ &\frac{1}{2} \cdot \ln((2\pi)^n |\Sigma|) = \frac{1}{2} \sum_i \sum_j \Sigma_{i,j}^{-1} \cdot E(x_i - a_i)(x_j - a_j) + \frac{1}{2} \cdot \ln((2\pi)^n |\Sigma|) = \frac{1}{2} \sum_i \sum_j \Sigma_{i,j}^{-1} \cdot \Sigma_{j,i} + \frac{1}{2} \cdot \\ \ln((2\pi)^n |\Sigma|) &= \frac{1}{2} \sum_i \sum_j \Sigma_{i,j}^{-1} \cdot \Sigma_{j,i} + \frac{1}{2} \cdot \ln((2\pi)^n |\Sigma|) = \frac{1}{2} \sum_i E_{i,i} + \frac{1}{2} \cdot \ln((2\pi)^n |\Sigma|) = \frac{1}{2} \cdot \ln((2\pi \cdot e)^n |\Sigma|) \end{aligned}$$