

Машинное обучение. Задание 2

Сопильняк Ольга

9 марта 2017 г.

1.1 Ответы в листьях регрессионного дерева

Пусть X_1, \dots, X_n — таргеты объектов обучающей выборки.

Рассмотрим первый случай. Пусть сначала $\bar{X} = \frac{X_1 + \dots + X_n}{n}$ — среднее значение таргета, тогда

$$EMSE(\bar{X}) = \frac{1}{n} \sum_{i=1}^n (\bar{X} - X_i)^2 = \frac{n\bar{X}^2 - 2\bar{X} \sum X_i + \sum X_i^2}{n}$$

Теперь пусть Y — таргет случайного объекта. Посчитаем матожидание:

$$EMSE(Y) = \frac{1}{n} E \sum_{i=1}^n (Y - X_i)^2 = \frac{nEY^2 - 2 \sum X_i EY + \sum X_i^2}{n}$$

Посмотрим на разность этих матожиданий и после некоторых преобразований получим:

$$EMSE(Y) - EMSE(\bar{X}) = \overline{X^2} - (\bar{X})^2 \geq 0$$

Поэтому среднее значение таргета лучше.

1.3 Unsupervised Decision Tree

$$f(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \cdot e^{-1/2(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

— многомерное нормальное распределение.

Его энтропия:

$$H(f) = - \int \cdots \int f(x) \ln f(x) dx$$

Распишем:

$$\int \cdots \int f(x) \left(\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) + \ln(2\pi)^{n/2} |\Sigma|^{1/2} \right) dx$$

$$\frac{1}{2} E \left(\sum (x_i - \mu_i) (\Sigma^{-1}) (x_j - \mu_j) \right) + \frac{1}{2} \ln((2\pi)^n |\Sigma|)$$

Первое слагаемое равно $\frac{n}{2}$, поэтому окончательно получаем:

$$\frac{1}{2} \ln((2\pi e)^n |\Sigma|)$$