

Машинное обучение.

Задание 1

Пашментов Никита

1 Задание 1

Покажите, что если в наивном байесовском классификаторе классы имеют одинаковые априорные вероятности, а плотность распределения признаков в каждом классе имеет вид $P(x^{(k)} | y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x^{(k)} - \mu_{yk})^2}{2\sigma^2}}$, $x^{(k)}, k = 1, \dots, n$ - признаки объекта x , то классификация сводится к отнесению объекта x к классу y , центр которого μ_y ближе всего к x .

1.1 Решение

$$\hat{y} = \arg \max_y \left(\prod_{i=1}^n P(x^{(i)} | y) P(y) \right) = \arg \max_y \left(\prod_{i=1}^n P(x^{(i)} | y) \right)$$

Будем искать максимум прологарифмированной функции:

$$\hat{y} = \arg \max_y \left(\sum_{i=1}^n \ln(P(x^{(i)} | y)) \right)$$

Подставим плотность из условия:

$$\hat{y} = \arg \max_y \left(-\frac{1}{2} \ln(2\pi\sigma^2) - \sum_{i=1}^n \frac{(x^{(i)} - \mu_{yi})^2}{2\sigma^2} \right) = \arg \min_y \sum_{i=1}^n (x^{(i)} - \mu_{yi})^2$$

Нетрудно заметить, что:

$$\hat{y} = \arg \min_y (\rho(x, \mu_y))$$

2 Задание 3

Утверждается, что метод одного ближайшего соседа асимптотически (при условии, что максимальное по всем точкам выборки расстояние до ближайшего соседа стремится к нулю) имеет матожидание ошибки не более чем вдвое больше по сравнению с оптимальным байесовским классификатором (который это матожидание минимизирует).

Покажите это, рассмотрев задачу бинарной классификации. Достаточно рассмотреть вероятность ошибки на фиксированном объекте x , т.к. математическое ожидание ошибок на выборке размера V будет просто произведением V на эту вероятность. Байесовский классификатор ошибается на объекте x с вероятностью:

$$E_B = \min\{P(1 | x), P(0 | x)\}$$

Условные вероятности будем считать непрерывными функциями от $x \in \mathbb{R}^m$, чтобы иметь возможность делать предельные переходы. Метод ближайшего соседа ошибается с вероятностью:

$$E_N = P(y \neq y_n)$$

Здесь y - настоящий класс x , а y_n - класс ближайшего соседа x_n к объекту x в предположении, что в обучающей выборке n объектов, равномерно заполняющих пространство.

Докажите исходное утверждение, выписав выражение для E_N (принадлежность к классам 0 и 1 для объектов x и x_n считать независимыми событиями) и осуществив предельный переход по n .

2.1 Решение

$$E_N = P(y_n = 1 | x_n)P(0 | x) + P(y_n = 0 | x_n)P(1 | x) \simeq 2P(1 | x)P(0 | x) \text{ (по непрерывности } P(y | x))$$

$$2P(1 | x)P(0 | x) \leq 2\min\{P(1 | x), P(0 | x)\}$$

Следовательно:

$$E_B \leq E_N \text{ - ч.т.д.}$$