

# Задачи по байесовскому подходу к классификации

Денисов Д.М.

Исходные данные:

$$\begin{aligned}p(x|y = -1) &= \frac{1}{\pi(1+x^2)} \sim C(0, 1), \\p(x|y = +1) &= \frac{1}{3}I_{[a,b]} \sim U(0, 3), \\ \lambda_- &= 2, \quad \lambda_+ = 1, \\p(y = -1) &= 0.4, \quad p(y = +1) = 0.6.\end{aligned}$$

## 1 Оптимальный байесовский классификатор

Имеем:

$$\begin{aligned}a^*(x) &= \operatorname{argmax}_y \lambda_y P(y)p(x|y) \\&= \operatorname{argmax}_y \{ \lambda_- P(y = -1)p(x|y = -1), \lambda_+ P(y = +1)p(x|y = +1) \} \\&= \operatorname{argmax}_y \left\{ \frac{0.8}{\pi(1+x^2)}, 0.2I_{[a,b]} \right\} \\&= \operatorname{argmax}_y \{ f(x), g(x) \}.\end{aligned}$$

Графики распределений  $f(x)$  и  $g(x)$  представлены на рис. 1.  
Из рис. 1 видно, что

$$g(x) \geq f(x) \iff x \in [x_0, 3].$$

Таким образом, оптимальный классификатор будет иметь вид

$$a^*(x) = \begin{cases} -1, & x \in (-\infty, x_0) \cup (3, +\infty) \\ +1, & x \in [x_0, 3] \end{cases}.$$

Координату  $x_0$  найдём следующим образом:

$$\begin{aligned}f(x_0) = g(x_0) = 0.2 &\implies \frac{0.8}{\pi(1+x_0^2)} = 0.2 \implies \\&\implies \pi(1+x_0^2) = 4 \implies x_0 = \sqrt{\frac{4}{\pi} - 1}.\end{aligned}$$

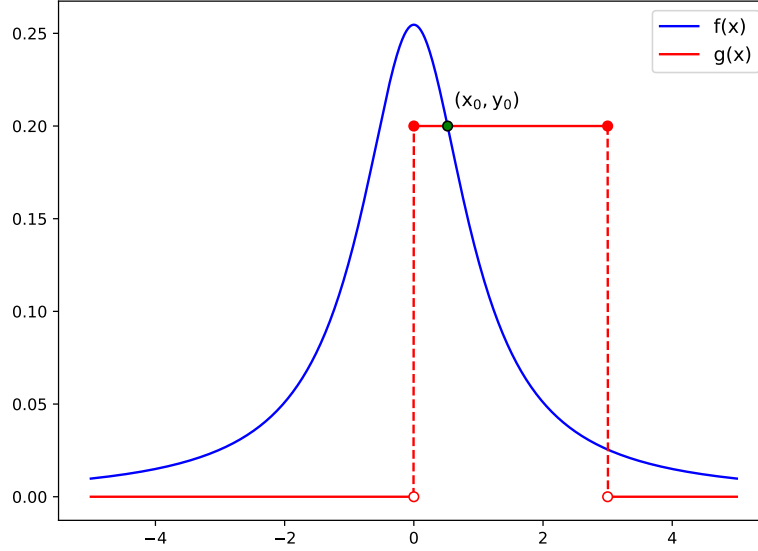


Рис. 1: Графики распределений

## 2 Оптимальный средний риск

Имеем:

$$\begin{aligned}
 R^* &= \iint L(a^*(x), y) p(x, y) dx dy \\
 &= \iint \lambda_y [a^*(x) \neq y] p(x|y) P(y) dx dy \\
 &= \iint \lambda_y (1 - [a^*(x) = y]) p(x|y) P(y) dx dy \\
 &= \int \min_y \lambda_y p(x|y) P(y) dx dy \\
 &= \int \min_y \{ \lambda_- p(x|y = -1) P(y = -1), \lambda_+ p(x|y = +1) P(y = +1) \} dx dy \\
 &= \int \min_y \{ f(x), g(x) \} dx dy.
 \end{aligned}$$

Аналогично пункту 1, получаем:

$$\begin{aligned}
 R^* &= \int_{-\infty}^{x_0} g(x) dx + \int_{x_0}^3 f(x) dx + \int_3^{+\infty} g(x) dx \\
 &= \int_0^{x_0} 0.2 dx + \int_{x_0}^3 \frac{0.8}{\pi(1+x^2)} dx \\
 &= 0.2x_0 + \frac{0.8}{\pi} (\arctan 3 - \arctan x_0) \approx 0.299958 \approx 0.3.
 \end{aligned}$$

### 3 Наивный байесовский классификатор

Имеем:

$$a^*(x_1, x_2) = \operatorname{argmax}_y P(y)p(x_1, x_2|y),$$

$$p(x_1, x_2|y) = p(x_1|y)p(x_2|y).$$

По исходным данным нетрудно определить априорные вероятности классов:

$$P(y = -1) = \frac{n_-}{n_+ + n_-} = 0.4,$$

$$P(y = +1) = \frac{n_+}{n_+ + n_-} = 0.6.$$

Предположим, что исходные данные соответствуют нормальному распределению, то есть

$$p(x_1|y = -1) \sim N(\mu_{11}, \sigma_{11}^2), \quad p(x_2|y = -1) \sim N(\mu_{21}, \sigma_{21}^2),$$

$$p(x_1|y = +1) \sim N(\mu_{12}, \sigma_{12}^2), \quad p(x_2|y = +1) \sim N(\mu_{22}, \sigma_{22}^2).$$

Определим параметры распределений по следующим формулам:

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i, \quad \sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}.$$

Получим следующие значения:

$$\begin{aligned} \mu_{11} &\approx -0.101, \quad \sigma_{11} \approx 1.17, \\ \mu_{21} &\approx 0.792, \quad \sigma_{21} \approx 2.014, \\ \mu_{12} &\approx 1.018, \quad \sigma_{12} \approx 1.011, \\ \mu_{22} &\approx 0.497, \quad \sigma_{22} \approx 1.918. \end{aligned}$$

Имеем:

$$\begin{aligned} a^*(x_1, x_2) &= \operatorname{argmax}_y P(y)p(x_1, x_2|y) \\ &= \operatorname{argmax}_y \{P(y = -1)p(x_1, x_2|y = -1), P(y = +1)p(x_1, x_2|y = +1)\} \\ &= \operatorname{argmax}_y \left\{ \begin{aligned} &P(y = -1)p(x_1|y = -1)p(x_2|y = -1), \\ &P(y = +1)p(x_1|y = +1)p(x_2|y = +1) \end{aligned} \right\} \\ &= \operatorname{argmax}_y \{f(x_1, x_2), g(x_1, x_2)\}. \end{aligned}$$

Очевидно, что

$$f(x) = \frac{0.4}{2\pi\sigma_{11}\sigma_{21}} e^{-\frac{1}{2} \left[ \left( \frac{x_1 - \mu_{11}}{\sigma_{11}} \right)^2 + \left( \frac{x_2 - \mu_{21}}{\sigma_{21}} \right)^2 \right]},$$

$$g(x) = \frac{0.6}{2\pi\sigma_{12}\sigma_{22}} e^{-\frac{1}{2} \left[ \left( \frac{x_1 - \mu_{12}}{\sigma_{12}} \right)^2 + \left( \frac{x_2 - \mu_{22}}{\sigma_{22}} \right)^2 \right]}.$$

Таким образом, оптимальный наивный классификатор будет иметь вид

$$a^*(x_1, x_2) = \begin{cases} -1, & f(x_1, x_2) \geq g(x_1, x_2) \\ +1, & f(x_1, x_2) < g(x_1, x_2) \end{cases}.$$