

# Resumo de Aprendizagem de Máquina 2014-2

Eduardo M. B. de A. Tenório

embat@cin.ufpe.br

CIn-UFPE

## Resumo

*Este documento tem por finalidade ser um resumo dos assuntos abordados na disciplina Aprendizagem de Máquina do período 2014-2 do CIn-UFPE, ministrada pelos professores Francisco Carvalho e Teresa Ludermit. A maioria do documento referencia o livro “Pattern Classification”, de Duda, Hart e Stork. Os códigos utilizados como exercício de fixação encontram-se em [github.com/embatbr/resumo-aprendizagem](https://github.com/embatbr/resumo-aprendizagem).*

## 1 Teoria da Decisão Bayesiana

Teoria da Decisão Bayesiana é uma abordagem estatística para a classificação de padrões, baseada em quantificar os tradeoffs associados a tomar uma determinada decisão (classificar) utilizando probabilidade e considerando os custos associados.

O **estado natural** é denotado por  $\omega$ , de modo que  $\omega = \omega_i$ , para  $i = 1, 2, \dots, c$ , significa que o exemplo foi classificado como pertencente à classe  $\omega_i$ . Cada uma dessas classes possui uma **probabilidade a priori**  $P(\omega_i)$ , com

$$\sum_{i=1}^c P(\omega_i) = 1, \quad (1)$$

refletindo o conhecimento prévio da chance de um elemento da classe  $\omega_i$  aparecer. A **regra**

**de decisão** seria: Decida  $\omega_i$  para  $\max_i P(\omega_i)$ . Neste caso a classe  $\omega_i$  sempre é escolhida e a probabilidade de erro é dada por:

$$P_{err}(\omega_i) = 1 - P(\omega_i). \quad (2)$$

Utilizando uma característica  $x$  que seja contínua e aleatória, sua **densidade de probabilidade estado-condicional** é dada por  $p(x|\omega)$ . Logo, a diferença entre  $p(x|\omega_i)$  e  $p(x|\omega_j)$  descreve a diferença da característica  $x$  entre as populações das classes  $\omega_i$  e  $\omega_j$ .

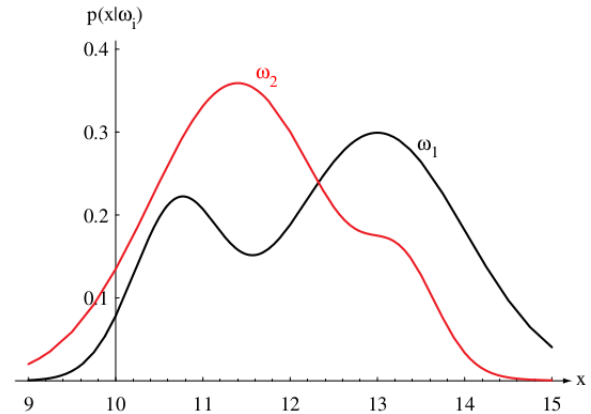


Figura 1: Para  $\omega_i = \omega_2$ , é mais frequente observar  $x$  entre 11 e 12 que  $x = 13$  (valor mais provável se  $\omega_i = \omega_1$ ).

Sabendo  $P(\omega_i)$  e  $p(x|\omega_i)$ , e medindo um valor  $x$ , a probabilidade conjunta de achar um padrão na classe  $\omega_i$  e com  $x$  é dado por:  $p(\omega_i, x) = P(\omega_i|x)p(x) = p(x|\omega_i)P(\omega_i)$ , que pela **fórmula de Bayes** fica:

$$P(\omega_i|x) = \frac{p(x|\omega_i)P(\omega_i)}{p(x)}, \quad (3)$$

com a evidência para  $c$  classes

$$p(x) = \sum_{j=1}^c p(x|\omega_j)P(\omega_j). \quad (4)$$

A fórmula de Bayes pode ser expressa em português como

$$posteriori = \frac{verossimilhanca \times priori}{evidencia}. \quad (5)$$

Fig. (2) mostra a probabilidade a posteriori das classes  $\omega_1$  e  $\omega_2$  para um conjunto de valores de  $x$ . A regra de decisão muda para: Decida  $\omega_i$  se  $\omega_i$  minimiza  $P(erro|x)$ , onde

$$P(erro|x) = \sum_{j \neq i} P(\omega_j|x), \quad (6)$$

ou simplesmente  $P(erro|x) = 1 - P(\omega_i|x)$ . Então a regra torna-se: Decida  $\omega_i$  para  $\max_i P(\omega_i|x)$

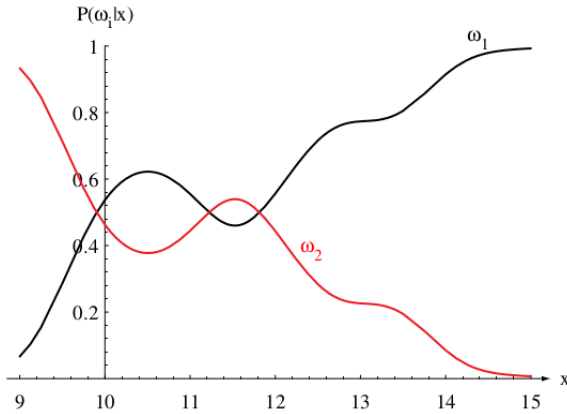


Figura 2: Probabilidades a posteriori para  $P(\omega_1) = \frac{2}{3}$  e  $P(\omega_2) = \frac{1}{3}$ , e para as densidades de probabilidade estado-condicional mostradas em Fig. (1).

Esta regra minimiza a probabilidade média de erro, dada por

$$P(erro) = \int_{-\infty}^{\infty} P(erro|x)p(x)dx. \quad (7)$$

É de fácil compreensão que a característica  $x$  pode ser trocada por um vetor de características  $\mathbf{x} = (x_1, x_2, \dots, x_d)$ , onde  $\mathbf{x}$  pertence ao espaço  $\mathbf{R}^d$ .

Outras ações além de apenas classificar um elemento podem ser tomadas, como por exemplo a **rejeição**: recusar-se a tomar uma decisão; uma opção válida quando o custo de ser indeciso não é tão alto. Para isso **funções de custo** são inseridas, permitindo tratar de situações onde alguns erros de classificação são mais importantes que outros.

Seja  $\{\omega_1, \dots, \omega_c\}$  o conjunto finito de  $c$  classes e seja  $\{\alpha_1, \dots, \alpha_a\}$  o conjunto finito de possíveis ações. A função de custo  $\lambda(\alpha_i|\omega_j)$  descreve o custo de tomar a ação  $\alpha_i$  quando a classe é  $\omega_j$ . Logo, observado um  $\mathbf{x}$  em particular, tomar a ação  $\alpha_i$  quando a classe é  $\omega_j$  leva a um custo esperado (**risco**)

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j)P(\omega_j|\mathbf{x}). \quad (8)$$

$R(\alpha_i|\mathbf{x})$  é chamado de **risco condicional**. Qualquer que seja o  $\mathbf{x}$  observado, o risco pode ser minimizado selecionando a ação que minimiza  $R(\alpha_i|\mathbf{x})$ .

A regra de decisão geral é uma função  $\alpha(\mathbf{x})$  que diz qual ação tomar para cada possível observação, ou seja, para cada  $\mathbf{x}$  a **função de decisão**  $\alpha(\mathbf{x})$  assume um dos  $a$  valores  $\alpha_1, \dots, \alpha_a$ . Logo, o risco global é dado por

$$R = \int R(\alpha(\mathbf{x}))p(\mathbf{x})d\mathbf{x}. \quad (9)$$

O risco global mínimo é chamado de **risco de Bayes**, denotado por  $R^*$ , sendo a melhor performance alcançável.