

### CT 720 Tópicos em Aprendizagem de Máquina e Classificação de Padrões



# 2-Teoria Bayesiana de Decisão

### Conteúdo

- 1. Introdução
- 2. Teoria Bayesiana de decisão: atributos contínuos
- 3. Classificação com taxa de erro mínima
- 4. Funções de discriminação e classificadores
- 5. Densidade normal
- 6. Funções de discriminação para densidade normal
- 7. Teoria Bayesiana de decisão: atributos discretos
- 8. Redes Bayesianas
- 9. Resumo

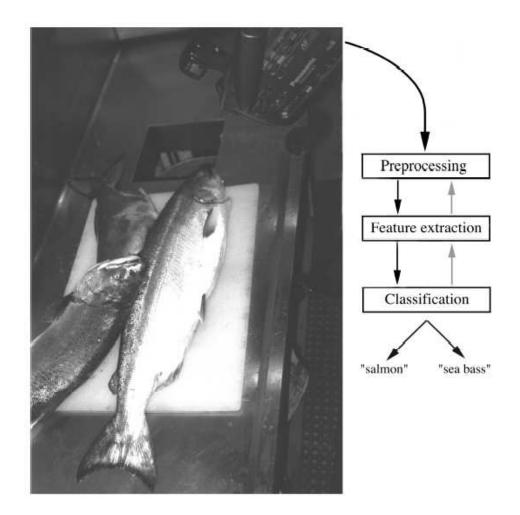
## 1-Introdução

#### Teoria Bayesiana de decisão

- abordagem estatística para reconhecimento padrões
- assume problema de decisão formulado probabilisticamente
- todos valores relevantes de probabilidade conhecidos
- identificação de sequências DNA

#### Este capítulo

- apresenta fundamentos da teoria
- teoria: formaliza procedimentos intuitivos



### Previsão próximo tipo peixe?

 $\omega$  = estado : variável aleatória

$$\omega = \omega_1$$
 sea bass

$$\omega = \omega_2$$
 salmon

#### Assumindo

 $P(\omega_1)$ : probabilidade *a priori* próximo tipo é *sea bass* 

 $P(\omega_1)$ : probabilidade *a priori* próximo tipo é *salmon* 

$$P(\omega_1) + P(\omega_2) = 1$$

#### Decidir tipo próximo peixe

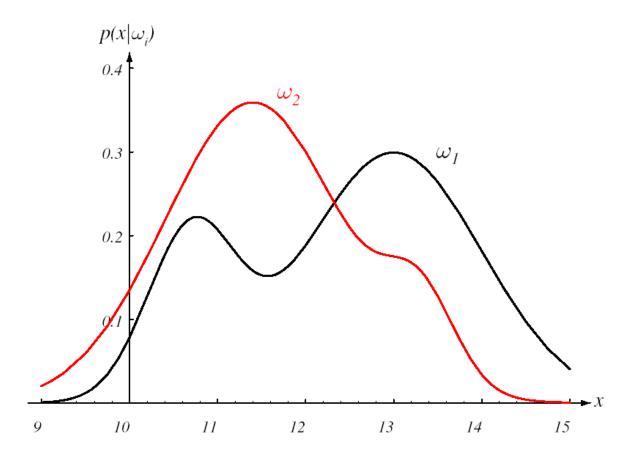
- classificação incorreta: mesmo custo
- informação disponíveis: somente  $P(\omega_1)$  e  $P(\omega_2)$

#### Regra de decisão

- $\square$  decidir  $\omega_1$  se  $P(\omega_1) > P(\omega_2)$ ; caso contrário decidir  $\omega_2$
- regra: mesma decisão, mesmo sabendo que existem 2 tipos
- desempenho depende da escolha de  $P(\omega_1)$  e  $P(\omega_2)$

#### Na prática

- temos mais informação
- e.g. medida luminosidade (x)
- densidade probabilidade condicional de classe  $p(x/\omega)$
- função densidade de probabilidade de x dado estado  $\omega$
- $-p(x/\omega_1)$  e  $p(x/\omega_2)$  descrevem diferenças de luminosidade



Supor que conhecemos:

- $P(\omega_1) e P(\omega_2)$
- $p(x|\omega_1) e p(x|\omega_2)$
- medida de luminosidade x
- Como *x* influencia a escolha correta do tipo?
  - estimativa correta do estado ?
  - como decidir o tipo (classe)

#### Classificação usando luminosidade como atributo

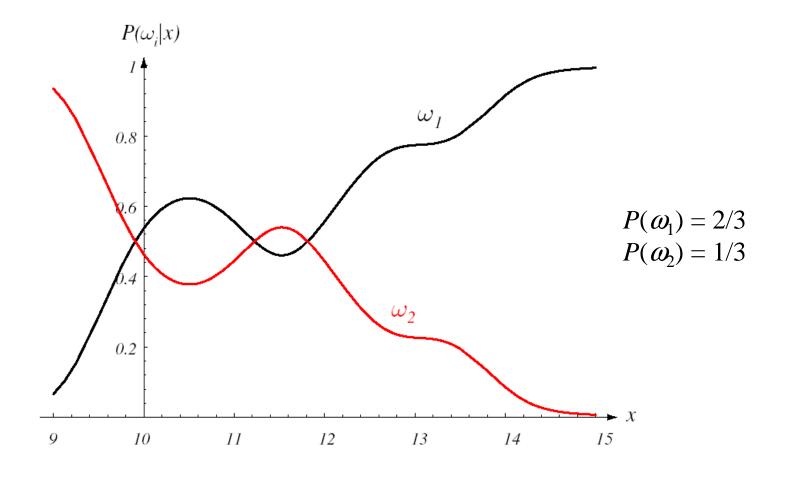
$$p(\omega_j, x) = P(\omega_j | x) p(x) = p(x | \omega_j) P(\omega_j)$$

$$P(\omega_j \mid x) = \frac{p(x \mid \omega_j)P(\omega_j)}{p(x)}$$
Bayes

$$P(\omega_j | x) = \frac{p(x | \omega_j) P(\omega_j)}{p(x)}$$
$$p(x) = \sum_{j=1}^{2} p(x | \omega_j) P(\omega_j)$$

$$P(causa / efeito) = \frac{p(efeito / causa)P(causa)}{P(efeito)}$$

$$posterior = \frac{likelyhood \times prior}{evidence}$$



### Regra de decisão de Bayes

 $\square$  decidir  $\omega_1$  se  $P(\omega_1|x) > P(\omega_2|x)$ ; caso contrário decidir  $\omega_2$  (\*)

Regra de Bayes minimiza a probabilidade de erro

$$P(error \mid x) = \begin{cases} P(\omega_1 \mid x) & \text{se decidimos } \omega_2 \\ P(\omega_2 \mid x) & \text{se decidimos } \omega_1 \end{cases}$$

- para uma observação *x* minimiza-se a probabilidade de erro:
  - $\square$  decidindo  $\omega_1$  se  $P(\omega_1|x) > P(\omega_2|x)$ ; ou  $\omega_2$  caso contrário
- esta regra minimiza a probabilidade de erro para qualquer x, em média?

– em média a probabilidade do erro é

$$P(error) = \int_{-\infty}^{\infty} P(error, x) dx = \int_{-\infty}^{\infty} P(error \mid x) p(x) dx$$

- $-\log a$  integral P(error) é mínima se P(error/x) é mínima, isto é, se
  - $\square$  decidir  $\omega_1$  se  $P(\omega_1|x) > P(\omega_2|x)$ ; caso contrário decidir  $\omega_2$
- $-P(error/x) = \min [P(\omega_1|x), P(\omega_2|x)]$
- Regra de decisão de Bayes (alternativa)
  - $\square$  decidir  $\omega_1$  se  $p(x|\omega_1)P(\omega_1) > p(x|\omega_2)P(\omega_2)$ ; caso contrário decidir  $\omega_2$

## 2-Teoria Bayesiana de decisão: atributos contínuos

#### Generalização

- uso de vários atributos
- mais de dois estados
- outras decisões e não só classificação
- critério mais geral que probabilidade de erro

#### Notação

 $\mathbf{x}$ : atributo,  $\mathbf{x} \in \mathbf{R}^d$ 

 $\mathbf{R}^d$ : espaço (Euclideano) de atributos

 $\{\omega_1, ..., \omega_c\}$ : conjunto (finito) de c estados (categorias)

 $\{\alpha_1, ..., \alpha_a\}$ : conjunto (finito) de *a* decisões (ações)

 $\lambda(\alpha_i, \omega_i)$ : loss function = custo decisão  $\alpha_i$  quando em  $\omega_i$ 

#### Regra de Bayes

$$P(\omega_j | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_j) P(\omega_j)}{p(\mathbf{x})}$$
$$p(\mathbf{x}) = \sum_{j=1}^{c} p(\mathbf{x} | \omega_j) P(\omega_j)$$

$$p(\mathbf{x}) = \sum_{j=1}^{c} p(\mathbf{x} \mid \omega_j) P(\omega_j)$$

- supor observação  $\mathbf{x}$  e ação  $\alpha_i$  correspondente
- se estado verdadeiro é  $\omega_i$ , então custo associado é  $\lambda(\alpha_i, \omega_i)$
- valor esperado do custo da ação  $\alpha_i$  será:

$$R(\alpha_i \mid \mathbf{x}) = \sum_{i=1}^{c} \lambda(\alpha_i \mid \omega_j) P(\omega_j \mid \mathbf{x})$$
 Risco condicional

- dado  $\mathbf{x}$ , que ação  $\alpha_i$  minimiza o risco condicional,  $\forall \mathbf{x}$ ?
- solução (ótima): regra de Bayes!

#### Regra (estratégia) de decisão

- $-\alpha(\mathbf{x})$ : fornece a decisão para cada valor de  $\mathbf{x}$
- para cada  $\mathbf{x}$ ,  $\alpha(\mathbf{x})$  assume um dos a valores  $\alpha_1$ , ...,  $\alpha_a$
- risco global: risco associado com uma estratégia de decisão

$$R = \int R(\alpha(\mathbf{x}) | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$
 Risco global

- se  $\alpha(\mathbf{x})$  é tal que o valor de  $R(\alpha_i(\mathbf{x}))$  é o menor possível  $\forall \mathbf{x}$ , então risco global é minimizado; isto motiva o seguinte:

#### Regra de Bayes

para minimizar risco global:

1 – calcular

$$R(\alpha_i \mid \mathbf{x}) = \sum_{i=1}^{c} \lambda(\alpha_i \mid \omega_j) P(\omega_j \mid \mathbf{x}), \quad i = 1, ..., a$$

2 – selecionar ação  $\alpha_i$  que minimiza  $R(\alpha_i|\mathbf{x})$ 

 $R^*$  = risco de Bayes

#### Exemplo com duas classes

- $-\alpha_1$ : decide que estado verdadeiro é  $\omega_1$
- $-\alpha_2$ : decide que estado verdadeiro é  $\omega_2$
- $-\lambda_{ij} = \lambda(\alpha_i, \omega_j)$  custo quando decisão  $\Rightarrow \omega_i$  mas verdadeiro é  $\omega_j$

risco condicional

$$R(\alpha_1 \mid \mathbf{x}) = \lambda_{11} P(\omega_1 \mid \mathbf{x}) + \lambda_{12} P(\omega_2 \mid \mathbf{x})$$
$$R(\alpha_2 \mid \mathbf{x}) = \lambda_{21} P(\omega_1 \mid \mathbf{x}) + \lambda_{22} P(\omega_2 \mid \mathbf{x})$$

 $\square$  decidir  $\omega_1$  se  $R(\alpha_1|\mathbf{x}) < R(\alpha_2|\mathbf{x})$ 

– alternativamente, em termos das probabilidades *a posteriori* 

$$\square$$
 decidir  $\omega_1$  se  $(\lambda_{21} - \lambda_{11})P(\omega_1 | \mathbf{x}) > (\lambda_{12} - \lambda_{22}) P(\omega_2 | \mathbf{x})$ 

- em geral  $\lambda_{21} > \lambda_{11}$  e  $\lambda_{12} > \lambda_{22}$
- utilizando Bayes, probabilidades a priori e densidades condicionais

$$\frac{P(\mathbf{x} \mid \omega_1)}{P(\mathbf{x} \mid \omega_2)} > \frac{(\lambda_{12} - \lambda_{22})}{(\lambda_{21} - \lambda_{11})} \frac{P(\omega_2)}{P(\omega_1)}$$

Razão de verosimilhança

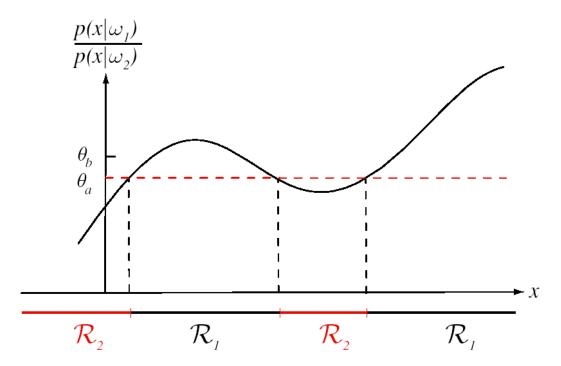
# 3-Classificação com taxa de erro mínima

$$\lambda(\alpha_i \mid \omega_j) = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases} \quad i, j = 1, ..., c$$

$$R(\alpha_i \mid \mathbf{x}) = \sum_{i=1}^{c} \lambda(\alpha_i \mid \omega_j) P(\omega_j \mid \mathbf{x})$$
$$= \sum_{i \neq j} P(\omega_j \mid \mathbf{x})$$
$$= 1 - P(\omega_i \mid \mathbf{x})$$

### Regra Bayes para taxa de erro mínima

- minimizar risco: selecionar ação que minimiza risco condicional
- minimizar taxa de erro de classificação:
  - $\Box \operatorname{decidir} \omega_i \operatorname{se} P(\omega_i | \mathbf{x}) > P(\omega_i | \mathbf{x}) \quad \forall i, j \quad (\text{ver *})$



# 4-Funções discriminação e classificadores

#### Função discriminação

$$g_i(\mathbf{x}), i = 1,..., c$$

 $\square$  classificador atribui classe  $\omega_i$  a **x** se  $g_i(\mathbf{x}) > g_i(\mathbf{x}) \quad \forall j \neq i$ 

#### Exemplos:

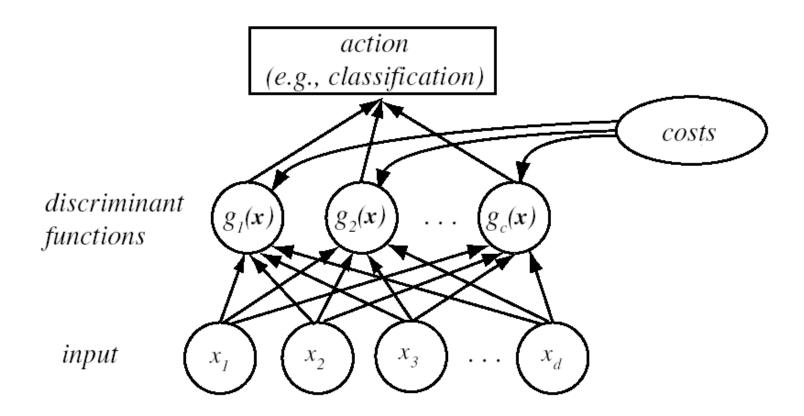
$$g_i(\mathbf{x}) = -R(\alpha_i|\mathbf{x})$$

risco condicional mínimo

$$g_i(\mathbf{x}) = P(\omega_i | \mathbf{x})$$

erro classificação mínimo

### Estrutura funcional de classificadores estatísticos



#### Propriedades

- funções discriminação não são únicas
- transformações:  $f(g_i(\mathbf{x}))$ , f monotônica crescente
- simplificação analítica e computacional
- exemplos de classificadores (erro classificação mínimo):

$$g_i(\mathbf{x}) = P(\omega_i \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \omega_i)P(\omega_i)}{\sum_{j=1}^{c} p(\mathbf{x} \mid \omega_j)P(\omega_j)}$$

$$g_i(\mathbf{x}) = p(\mathbf{x} \mid \omega_i) P(\omega_i)$$

$$g_i(\mathbf{x}) = \ln p(\mathbf{x} \mid \omega_i) + \ln P(\omega_i)$$

- formas funções diferentes, mas regras de decisão equivalentes
- efeito: dividir o espaço de atributos em c regiões distintas
- $\square$  se  $g_i(\mathbf{x}) > g_i(\mathbf{x}) \quad \forall j \neq i \text{ então } \mathbf{x} \in \mathcal{R}_i$
- $-\mathcal{R}_i \cap \mathcal{R}_j = \emptyset$
- $-\mathcal{R}_1,...,\mathcal{R}_c$  formam uma partição do espaço de atributos

– exemplo: duas categorias (classes)

$$\square$$
 atribuir  $\omega_1$  a **x** se  $g_1(\mathbf{x}) > g_2(\mathbf{x}) \quad \forall j \neq i$ 

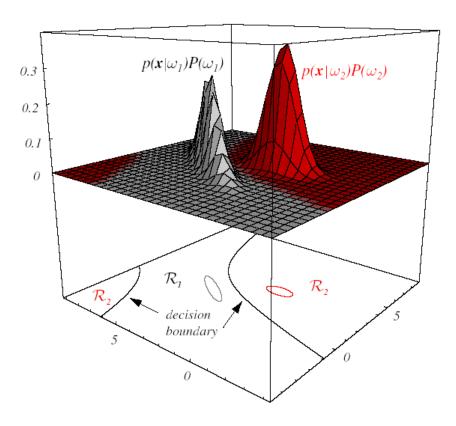
$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

$$\square$$
 atribuir  $\omega_1$  a **x** se  $g(\mathbf{x}) > 0$ 

$$g(\mathbf{x}) = P(\omega_1 \mid \mathbf{x}) - P(\omega_2 \mid \mathbf{x})$$

$$g(\mathbf{x}) = \ln \frac{p(\mathbf{x} \mid \omega_1)}{p(\mathbf{x} \mid \omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$
(\*\*)

(\*\*\*)



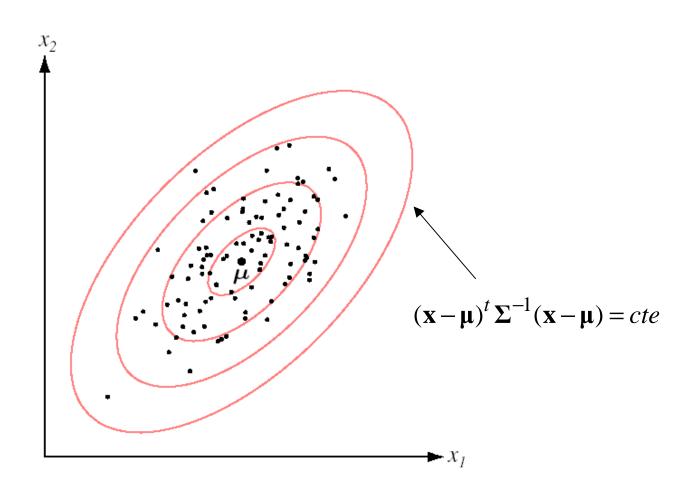
## 5-Densidade normal

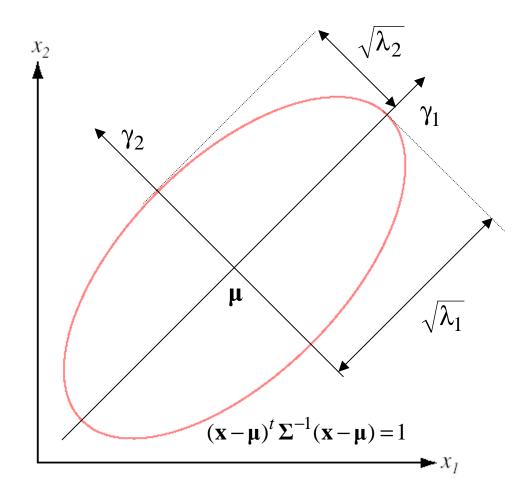
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\mathbf{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

$$\boldsymbol{\mu} = E[\mathbf{x}] = \int_{-\infty}^{\infty} \mathbf{x} \, p(\mathbf{x}) d\mathbf{x}, \quad \boldsymbol{\mu}_i = E[x_i]$$

$$\Sigma = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t] = \int_{-\infty}^{\infty} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t p(\mathbf{x}) d\mathbf{x}$$

$$\Sigma = [\sigma_{ij}], \ \sigma_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)], \ \Sigma > 0$$





#### Transformações lineares

– combinações lineares de variáveis aleatórias normais são normais

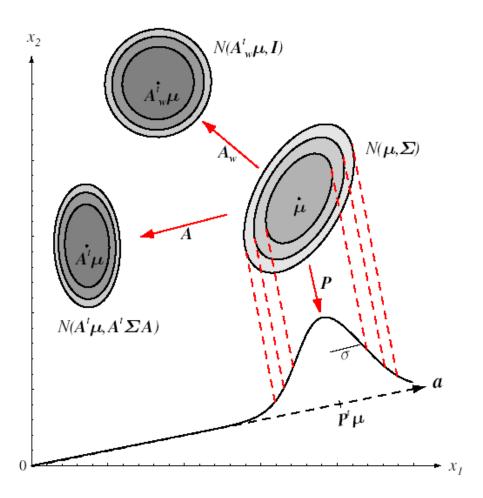
$$p(\mathbf{x}) = N(\mathbf{\mu}, \mathbf{\Sigma})$$

$$\mathbf{A} (d \times k)$$

$$\mathbf{y} = \mathbf{A}^{t}\mathbf{x} (k \times 1)$$

$$p(\mathbf{y}) = N(\mathbf{A}^{t}\mathbf{\mu}, \mathbf{A}^{t}\mathbf{\Sigma}\mathbf{A})$$

- se  $\mathbf{A} = \mathbf{a} (d \times 1)$ , k = 1,  $\|\mathbf{a}\| = 1$  então  $y = \mathbf{a}^t \mathbf{x}$  (escalar que representa projeção de  $\mathbf{x}$  ao longo de  $\mathbf{a}$ )  $\mathbf{a}^t \sum \mathbf{a}$  variância da projeção de  $\mathbf{x}$  ao longo de  $\mathbf{a}$ 



#### Observações

$$r^2 = (\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

#### Distância de Mahalanobis

$$V = V_d \mid \mathbf{\Sigma} \mid^{1/2} r^d$$

Volume hiperelipsóide

$$V_d = \begin{cases} \pi^{d/2} / (d/2)! & d \ par \\ 2^d \pi^{(d-1)/2} (d-1/2)! / d! & d \ impar \end{cases}$$

Volume hiperesfera

# 6-Funções discriminação p/ densidade normal

$$g_i(\mathbf{x}) = \ln p(\mathbf{x} \mid \omega_i) + \ln P(\omega_i)$$

$$p(\mathbf{x} | \omega_i) \sim N(\mathbf{\mu}_i, \mathbf{\Sigma}_i)$$

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i)$$

• Caso 1:  $\sum_{i} = \sigma^{2} \mathbf{I}$ 

$$g_i(\mathbf{x}) = -\frac{\|\mathbf{x} - \mathbf{\mu}_i\|^2}{2\sigma^2} + \ln P(\omega_i)$$

$$g_i(\mathbf{x}) = -\frac{1}{2\sigma^2} [\mathbf{x}^t \mathbf{x} + 2\mathbf{\mu}^t \mathbf{x} + \mathbf{\mu}_i^t \mathbf{\mu}_i] + \ln P(\omega_i)$$

termo quadrático é o mesmo  $\forall i$ 

$$g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0}$$
 Máquina linear

$$\mathbf{w}_{i} = \frac{1}{\sigma^{2}} \mathbf{\mu}_{i}, \qquad w_{io} = \frac{-1}{2\sigma^{2}} \mathbf{\mu}_{i}^{t} \mathbf{\mu}_{i} + \ln P(\omega_{i})$$

Limiar (bias) para i-ésima classe

$$g_i(\mathbf{x}) - g_j(\mathbf{x}) = 0$$

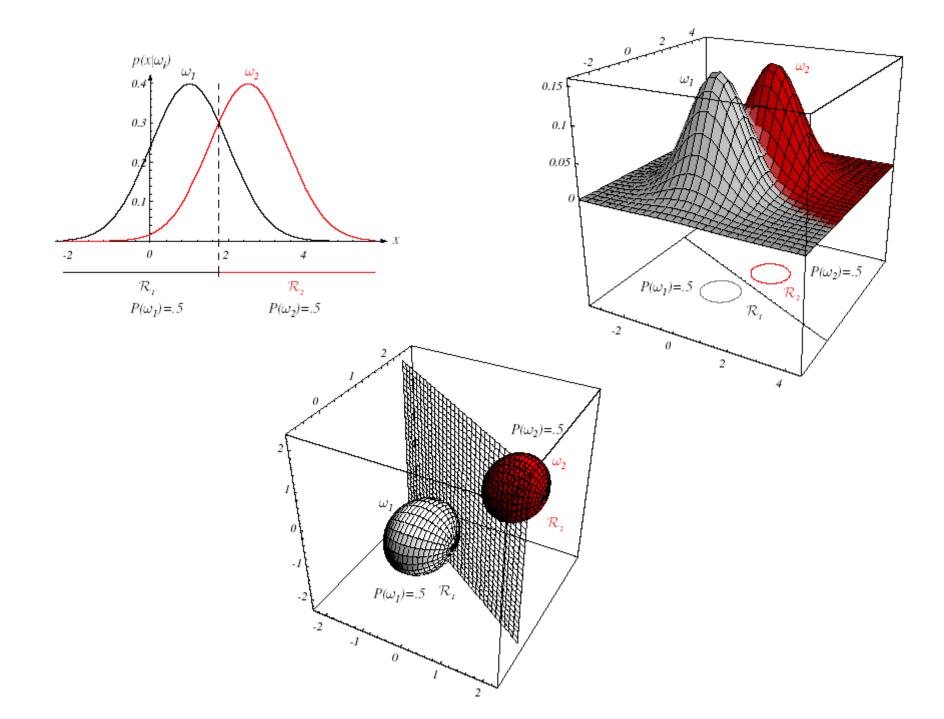
$$\mathbf{w}^t(\mathbf{x} - \mathbf{x}_0) = 0$$

$$\mathbf{w} = \mathbf{\mu}_i - \mathbf{\mu}_j$$

Hiperplano separa  $\mathcal{R}_i$  e  $\mathcal{R}_j$  passa por  $\mathbf{x}_o$  e é ortogonal à reta que une as médias

$$\mathbf{x}_{o} = \frac{1}{2} (\mathbf{\mu}_{i} + \mathbf{\mu}_{j}) - \frac{\sigma^{2}}{\|\mathbf{\mu}_{i} - \mathbf{\mu}_{j}\|^{2}} \ln \frac{P(\omega_{i})}{P(\omega_{j})} (\mathbf{\mu}_{i} - \mathbf{\mu}_{j})$$

$$P(\omega_i) = P(\omega_j) \implies \mathbf{x}_o = \frac{1}{2}(\mathbf{\mu}_i + \mathbf{\mu}_j)$$

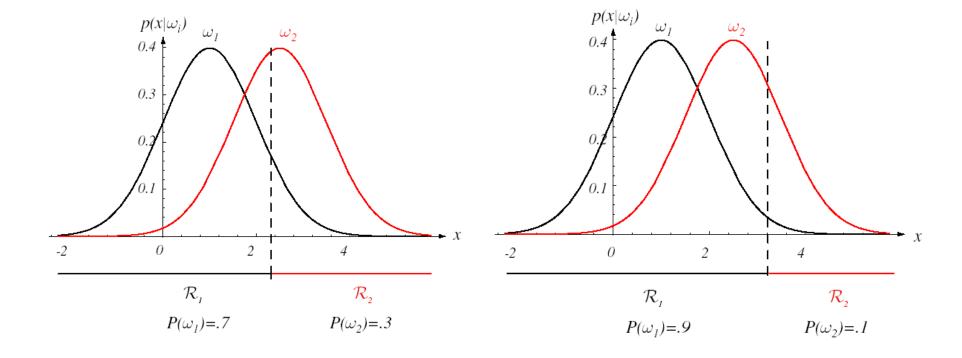


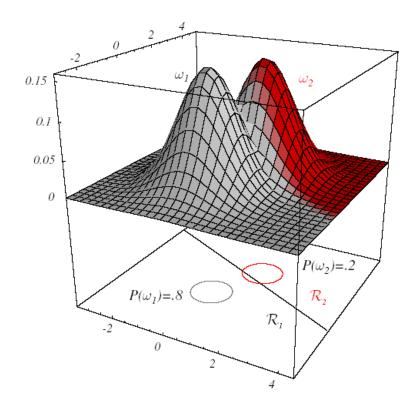
$$P(\omega_i) \neq P(\omega_i)$$

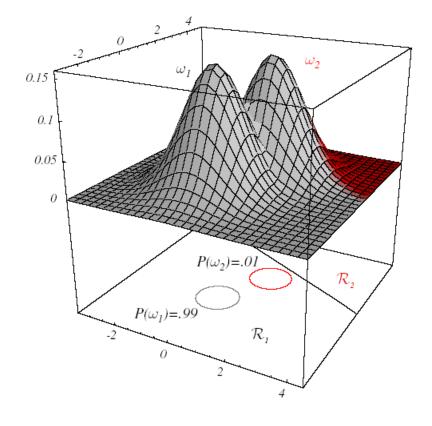
 $\mathbf{x}_{o}$  se afasta da média mais provável

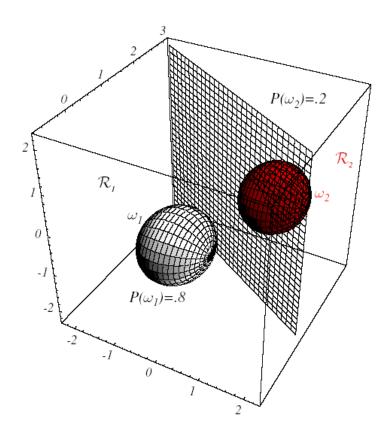
$$P(\omega_i) = P(\omega_j), \quad \forall i, j \quad \text{termo} \quad \ln P(\omega_i) \text{ irrelevante}$$

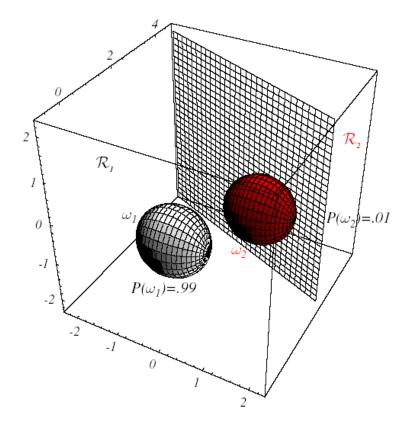
$$g_i(\mathbf{x}) = \|\mathbf{x} - \mathbf{\mu}_i\|^2$$
 Classificador distância mínima











• Caso 2:  $\sum_{i} = \sum_{i}$ 

independente de i

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{\mu})^t \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{\mu}) + \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\mathbf{\Sigma}_i| + \ln P(\omega_i)$$

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{\mu})^t \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{\mu}) + \ln P(\omega_i)$$

se  $P(\omega_i) = P(\omega_i)$ ,  $\forall i, j$  termo  $\ln P(\omega_i)$  irrelevante, logo:

$$g_i(\mathbf{x}) = \frac{1}{2} (\mathbf{x} - \mathbf{\mu})^t \mathbf{\Sigma}^{-1} (\mathbf{x} - \mathbf{\mu})$$
 Classificador distância (Mahalanobis) mínima

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{\mu})^t \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{\mu}) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\mathbf{\Sigma}_i| + \ln P(\omega_i)$$

eliminando  $\mathbf{x}^t \sum^{-1} \mathbf{x}$  da expansão de  $(\mathbf{x} - \boldsymbol{\mu})^t \sum^{-1} (\mathbf{x} - \boldsymbol{\mu})$  (independe de i)

$$g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{io}$$
 Máquina linear

$$\mathbf{w}_i = \mathbf{\Sigma}^{-1} \mathbf{\mu}_i \qquad w_{io} = -\frac{1}{2} \mathbf{\mu}_i^t \mathbf{\Sigma}^{-1} \mathbf{\mu}_i + \ln P(\omega_i)$$

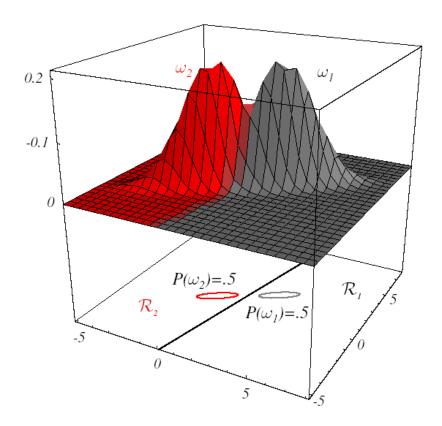
se  $\mathcal{R}_i$ e  $\mathcal{R}_j$  são contíguas, a superfície de decisão é um hiperplano

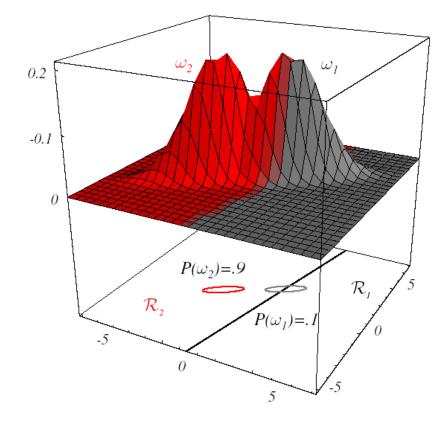
$$\mathbf{w}^t(\mathbf{x} - \mathbf{x}_0) = 0$$

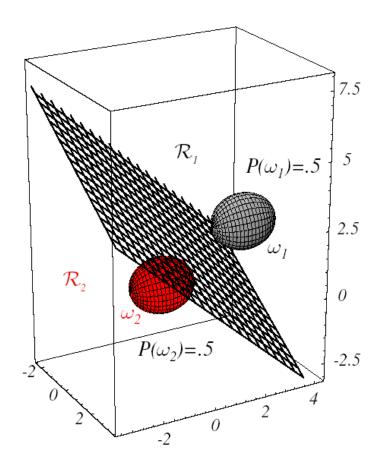
$$\mathbf{w} = \mathbf{\Sigma}^{-1}(\mathbf{\mu}_i - \mathbf{\mu}_j)$$

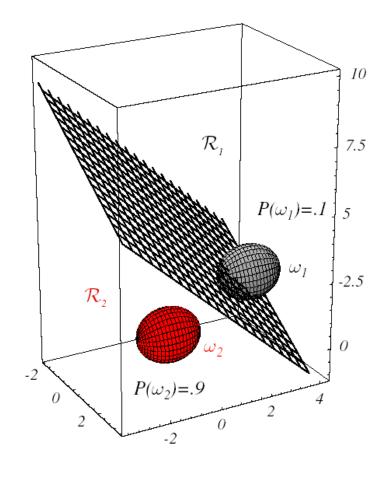
$$\mathbf{x}_{o} = \frac{1}{2} (\mathbf{\mu}_{i} + \mathbf{\mu}_{j}) - \frac{\ln[P(\omega_{i})/P(\omega_{j})]}{(\mathbf{\mu}_{i} - \mathbf{\mu}_{j})^{t} \Sigma^{-1} (\mathbf{\mu}_{i} - \mathbf{\mu}_{j})} (\mathbf{\mu}_{i} - \mathbf{\mu}_{j})$$

- hiperplano separa  $\mathcal{R}_i$  e  $\mathcal{R}_j$  não é ortogonal à reta que une as médias
- hiperplano intercepta esta reta em  $\mathbf{x}_{0}$
- $-\operatorname{se} P(\omega_i) = P(\omega_i), \ \forall i, j \ \operatorname{então} \operatorname{está} \operatorname{no} \operatorname{meio} \operatorname{das} \operatorname{médias}$
- se  $P(\omega_i)$  ≠  $P(\omega_j)$ , então hiperplano se afasta da média mais provável









• Caso 3:  $\sum_{i}$  = arbitrária

único termo independente de i

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{\mu})^t \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{\mu}) + \frac{d}{2} \ln 2\pi + \frac{1}{2} \ln |\mathbf{\Sigma}_i| + \ln P(\omega_i)$$

expandindo

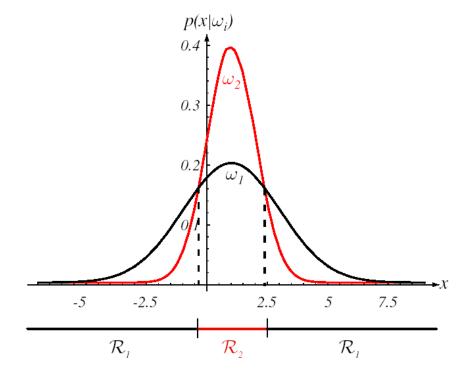
$$g_i(\mathbf{x}) = \mathbf{x}^t \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^t \mathbf{x} + w_{io}$$

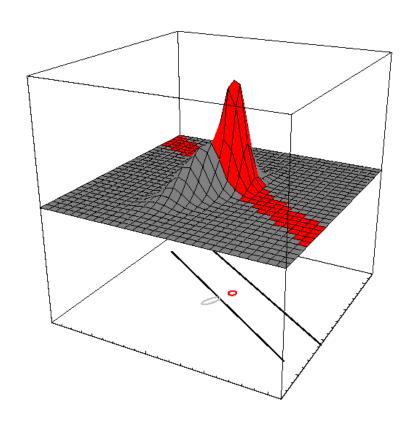
Classificador quadrático

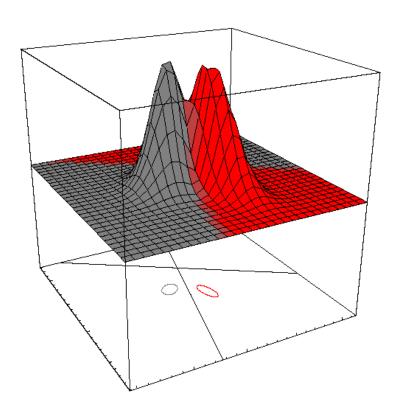
$$\mathbf{W}_i = -\frac{1}{2} \mathbf{\Sigma}_i^{-1} \qquad \mathbf{w}_i = \mathbf{\Sigma}_i^{-1} \mathbf{\mu}_i$$

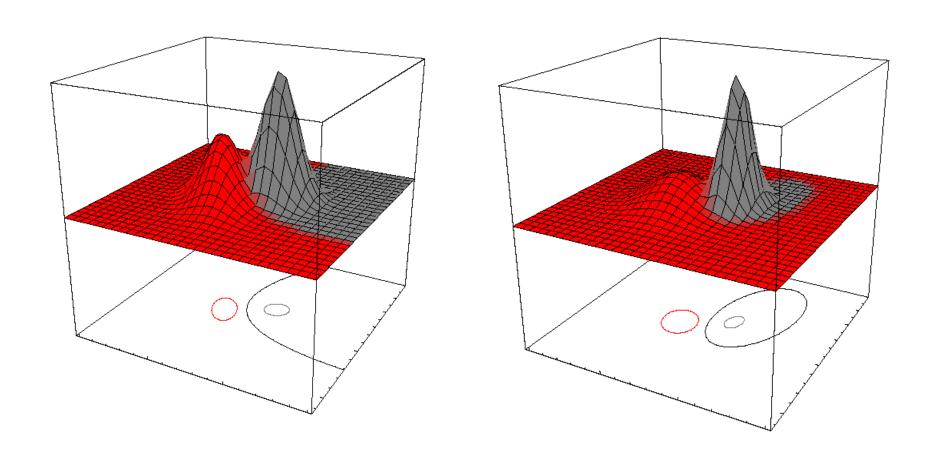
$$w_{io} = -\frac{1}{2} \boldsymbol{\mu}_i^t \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i^{-1}| + \ln P(\omega_i)$$

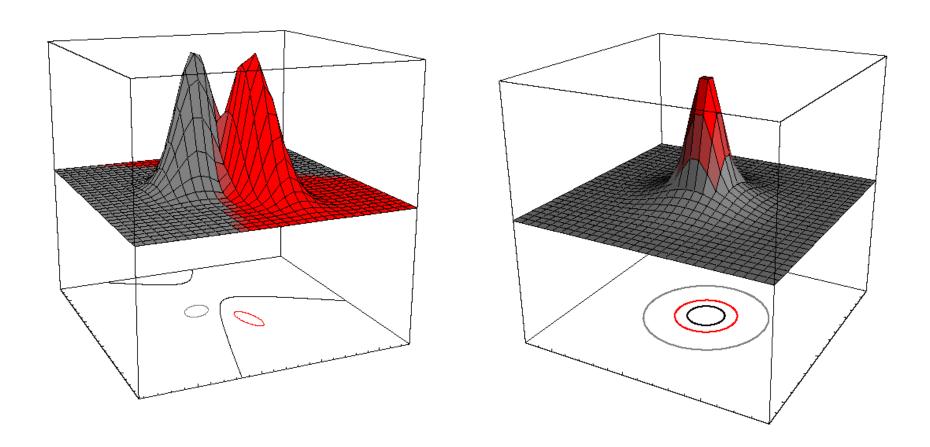
- caso com duas classes
- superfícies de decisão são hiperquadráticas
  - hiperplanos
  - hiperesferas
  - hiperelipsóides
  - hiperparabolóides
- regiões (decisão) não necessariamente conectadas

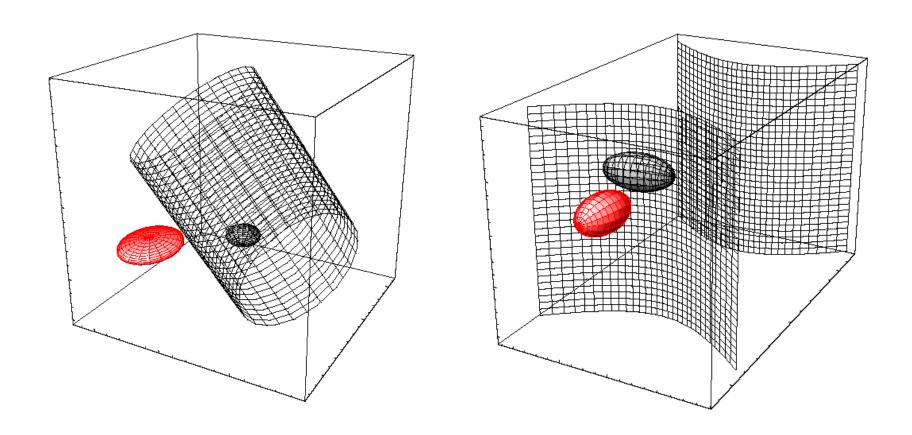


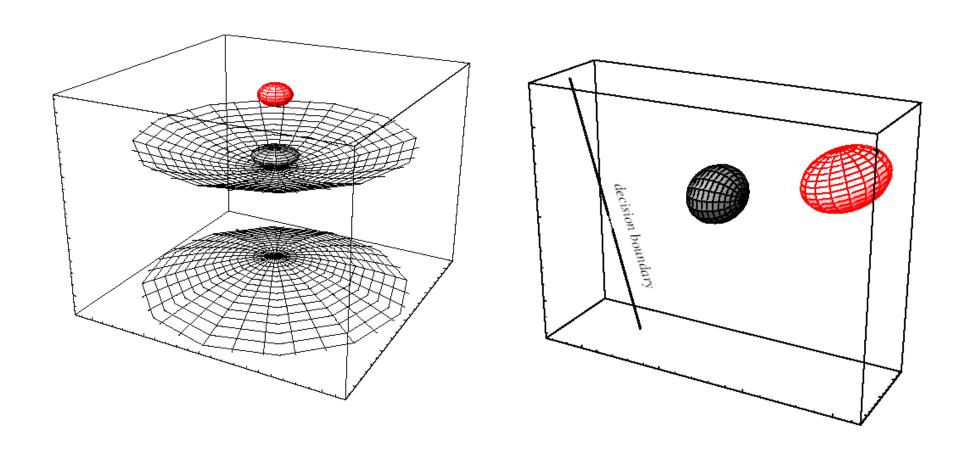




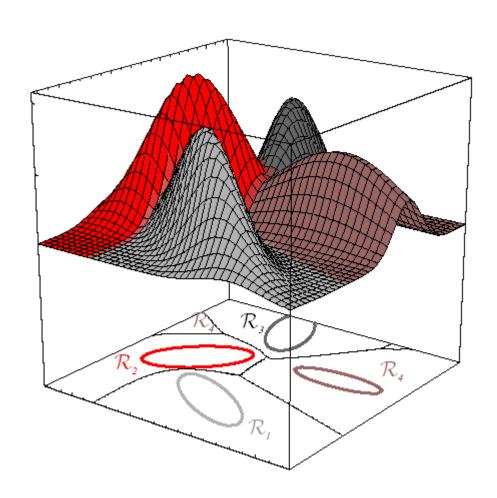




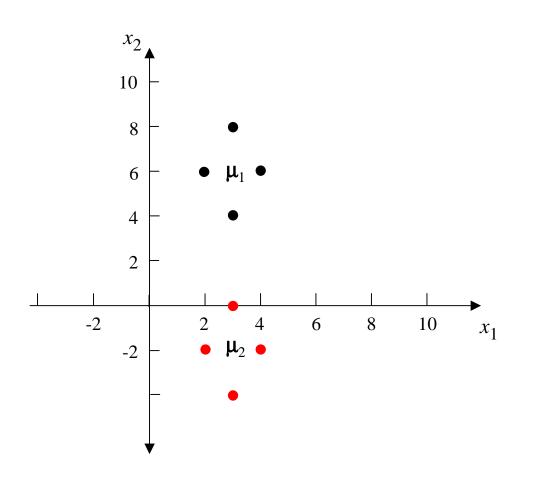




## Quatro classes



Exemplo: região de decisão, dados Gaussianos



$$\boldsymbol{\mu}_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix} \ \boldsymbol{\Sigma}_1 = \begin{bmatrix} 1/2 & 0 \\ 0 & 2 \end{bmatrix}$$

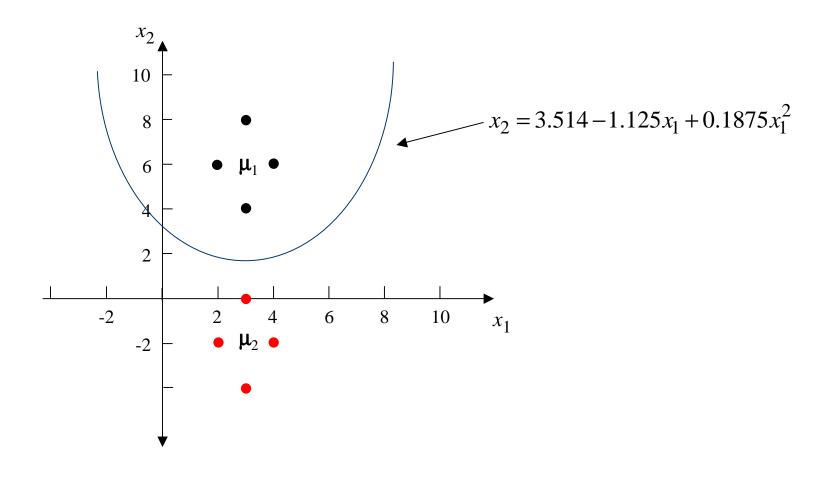
$$\boldsymbol{\mu}_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix} \quad \boldsymbol{\Sigma}_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$P(\omega_1) = P(\omega_2) = 0.5$$

$$g_1(\mathbf{x}) = \mathbf{x}^t \mathbf{W}_1 \mathbf{x} + \mathbf{w}_1^t \mathbf{x} + w_{10}$$

$$g_1(\mathbf{x}) = \mathbf{x}^t \mathbf{W}_1 \mathbf{x} + \mathbf{w}_1^t \mathbf{x} + w_{10}$$

$$g_1(\mathbf{x}) = g_2(\mathbf{x})$$



# 7-Teoria Bayesiana de decisão: atributos discretos

$$\mathbf{x} \in \{\mathbf{v}_1, ...., \mathbf{v}_m\}$$

$$P(\omega_j \mid \mathbf{x}) = \frac{P(\mathbf{x} \mid \omega_j)P(\omega_j)}{P(\mathbf{x})}$$

Regra de Bayes

$$P(\mathbf{x}) = \sum_{j=1}^{c} P(\mathbf{x} \mid \omega_j) P(\omega_j)$$

$$\alpha^* = \arg\min_i R(\alpha_i \mid \mathbf{x})$$

Regra de decisão mínimo risco

Regra de Bayes para taxa de erro mínima

Funções de discriminação

$$g_{i}(\mathbf{x}) = P(\omega_{i} \mid \mathbf{x}) = \frac{P(\mathbf{x} \mid \omega_{i})P(\omega_{i})}{\sum_{j=1}^{c} P(\mathbf{x} \mid \omega_{j})P(\omega_{j})}$$
$$g_{i}(\mathbf{x}) = P(\mathbf{x} \mid \omega_{i})P(\omega_{i})$$

$$g_i(\mathbf{x}) = \ln P(\mathbf{x} \mid \omega_i) + \ln P(\omega_i)$$

### Exemplo: atributos binários independentes

- duas categorias (classes)
- cada componente é um valor binário
- componentes s\(\tilde{a}\) independentes

$$\mathbf{x} = (x_1, \dots, x_d)^t, \ x_i \in \{0,1\}$$

$$p_i = \Pr[x_i = 1 \mid \omega_1)$$

$$q_i = \Pr[x_i = 1 \mid \omega_2)$$

$$P(\mathbf{x} \mid \omega_1) = \prod_{i=1}^{d} p_i^{x_i} (1 - p_i)^{1 - x_i}$$

$$P(\mathbf{x} \mid \omega_2) = \prod_{i=1}^{d} q_i^{x_i} (1 - q_i)^{1 - x_i}$$

### - Razão de verosimilhança

$$\frac{P(\mathbf{x} \mid \omega_1)}{P(\mathbf{x} \mid \omega_2)} = \prod_{i=1}^{d} \left(\frac{p_i}{q_i}\right)^{x_i} \left(\frac{1 - p_i}{1 - q_i}\right)^{1 - x_i}$$

$$g(\mathbf{x}) = \ln \frac{P(\mathbf{x} \mid \omega_1)}{P(\mathbf{x} \mid \omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$
 (\*\*)

$$g(\mathbf{x}) = \sum_{i=1}^{d} \left[ x_i \ln \frac{p_i}{q_i} + (1 - x_i) \ln \frac{1 - p_i}{1 - q_i} \right] + \ln \frac{P(\omega_1)}{P(\omega_2)}$$
 Linear em  $x_i$ 

$$g(\mathbf{x}) = \sum_{i=1}^{d} w_i x_i + w_0$$

$$w_i = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)}, \quad i = 1,...,d$$

$$w_0 = \sum_{i=1}^{d} \ln \frac{(1 - q_i)}{(1 - p_i)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

lembrar: decidir  $\omega_1$  se  $g(\mathbf{x}) > 0$  e  $\omega_2$  se  $\mathbf{g}(\mathbf{x}) \ge 0$  (\*\*\*)

#### Análise

- $-g(\mathbf{x})$  é uma combinação linear (ponderada) das componentes de  $\mathbf{x}$
- valor de  $w_i$  é a relevância de uma resposta  $x_i = 1$  na classificação
- $-\operatorname{se} p_i = q_i$  então valor de  $x_i$  é irrelevante ( $w_i = 0$ )
- $-\operatorname{se} p_i > q_i \operatorname{então} (1 p_i) < (1 q_i), w_i > 0 \operatorname{e} x_i = 1 \Rightarrow w_i \operatorname{votos} \operatorname{para} \omega_1$
- $-\operatorname{se} p_i < q_i \operatorname{então} (1 p_i) > (1 q_i), w_i < 0 \operatorname{e} x_i = 1 \Longrightarrow |w_i| \operatorname{votos} \operatorname{para} \omega_2$

### Exemplo: dados binários 3-d

- duas categorias (classes)
- cada componente é um valor binário
- componentes s\(\tilde{a}\) independentes

$$P(\omega_1) = P(\omega_2) = 0.5$$

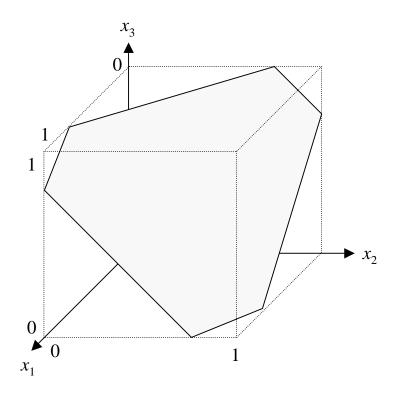
$$p_i = 0.5, q_i = 0.8, i = 1,2,3$$

$$w_i = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)}, \quad i = 1,...,d$$

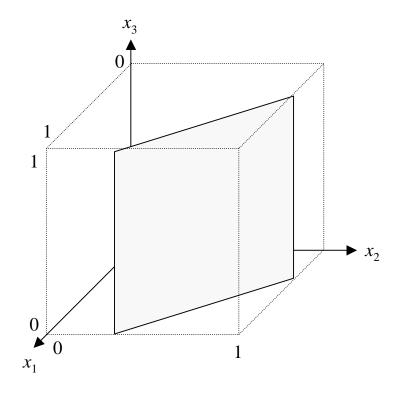
$$w_0 = \sum_{i=1}^{d} \ln \frac{(1-q_i)}{(1-p_i)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

$$w_i = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)}, \quad i = 1,...,d$$
  $w_i = \ln \frac{0.8(1-0.5)}{0.5(1-0.8)} = 1.3863 \quad i = 1,...,3$ 

$$w_0 = \sum_{i=1}^{3} \ln \frac{(1-0.8)}{(1-0.5)} + \ln \frac{0.5}{0.5} = -1.75$$



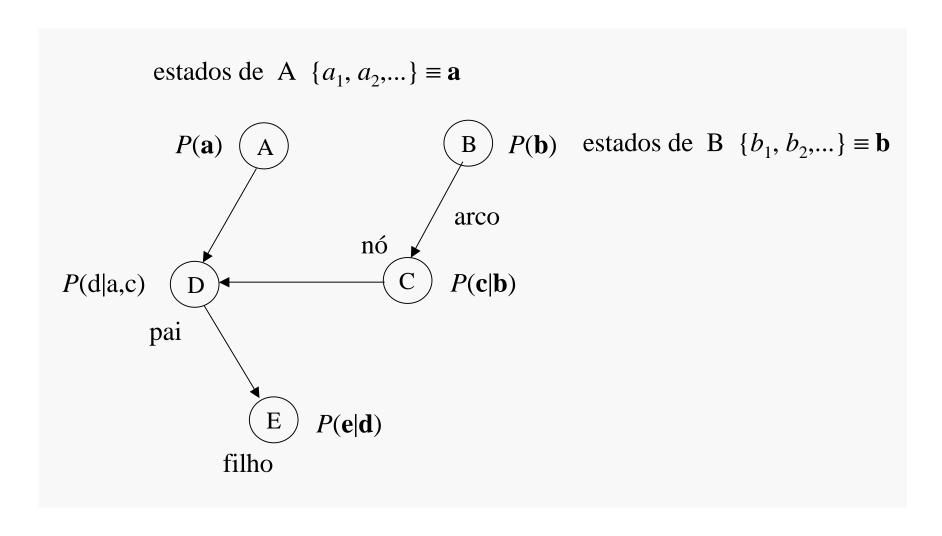
$$p_i = 0.8$$
  $q_i = 0.5$   $i = 1,2,3$ 

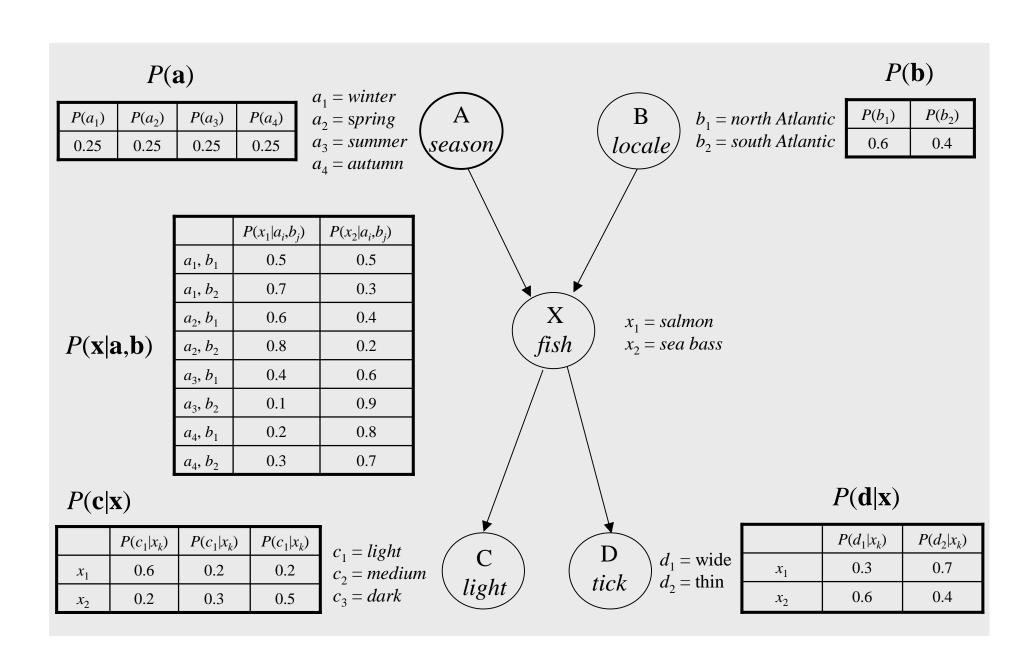


$$p_i = 0.8$$
  $q_i = 0.5$ ,  $i = 1,2$   
 $p_3 = q_3$ 

# 8-Redes Bayesianas

- Conhecimento sobre distribuições
  - parâmetros de distribuições
  - dependência/independência estatística
  - relações causais entre variáveis
- Redes Bayesianas
  - explora informação estrutural no raciocínio com variáveis
  - usa relações probabilísticas entre variáveis
  - assume relações causais

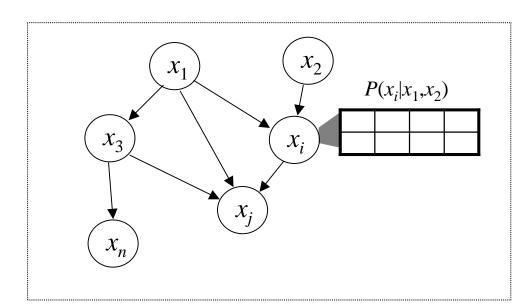




 $P(a_3,b_1,x_2,c_3,d_2) = P(a_3)P(b_1)P(x_2 \mid a_3,b_1)P(c_3 \mid x_2)P(d_2 \mid x_2) = 0.25 \times 0.6 \times 0.4 \times 0.5 \times 0.4 = 0.012$ 

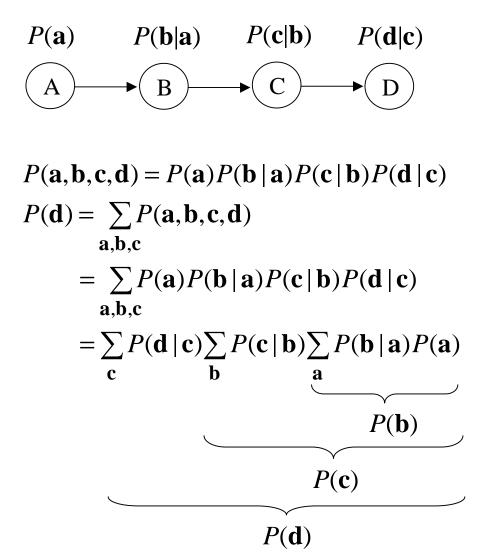
### Redes Bayesianas formalmente

- grafo acíclico
- nó: variável aleatória (atributo)
- arco: efeito, causa (A afeta B  $\rightarrow$  B condicionado a A)
- cada nó condicionalmente independente dos não descendentes
- representa probabilidade conjunta das variáveis



$$P(x_1,...,x_n) = \prod_{i=1}^n P(x_i \mid PaiDe(x_i))$$

### Exemplo



em geral: dados os valores de algumas variáveis (evidência: e)
 qual é o valor de uma configuração das outras variáveis (x)?

$$P(\mathbf{x} \mid \mathbf{e}) = \frac{P(\mathbf{x}, \mathbf{e})}{P(\mathbf{e})} = \alpha P(\mathbf{x}, \mathbf{e})$$

- Exemplo: salmon and sea bass
  - probabilidade peixe veio do Atlântico Norte  $(b_1)$
  - sabendo que é primavera (spring  $a_2$ )
  - peixe é salmão (salmon  $x_1$ ) claro (light  $c_1$ )

$$P(b_1 | a_2, x_1, c_1)$$
?

- Em classificação (erro mínimo): salmon or sea bass?
  - sabe-se que
    - peixe é claro  $(c_1)$
    - origem é Atlântico Norte  $(b_2)$
  - não se sabe:
    - estação do ano
    - espessura
  - problema de classificação:

$$P(x_1 | c_1, b_2)$$
?

$$P(x_2 | c_1, b_2)$$
?

$$P(x_{1} | c_{1}, b_{2}) = \frac{P(x_{1}, c_{1}, b_{2})}{P(c_{1}, b_{2})} = \alpha \sum_{\mathbf{a}, \mathbf{d}} P(x_{1}, \mathbf{a}, b_{2}, c_{1}, \mathbf{d})$$

$$= \alpha \sum_{\mathbf{a}, \mathbf{d}} P(\mathbf{a}) P(b_{2}) P(x_{1} | \mathbf{a}, b_{2}) P(c_{1} | x_{1}) P(\mathbf{d} | x_{1})$$

$$= \alpha P(b_{2}) P(c_{1} | x_{1}) \left[ \sum_{\mathbf{a}} P(\mathbf{a}) P(x_{1} | \mathbf{a}, b_{2}) \right] \left[ \sum_{\mathbf{d}} P(\mathbf{d} | x_{1}) \right]$$

$$= \alpha P(b_{2}) P(c_{1} | x_{1})$$

$$\times [P(a_{1}) P(x_{1} | a_{1}, b_{2}) + P(a_{2}) P(x_{1} | a_{2}, b_{2})$$

$$+ P(a_{3}) P(x_{1} | a_{3}, b_{2}) + P(a_{4}) P(x_{1} | a_{4}, b_{2})]$$

$$\times [P(d_{1} | x_{1}) + P(d_{2} | x_{1})]$$

$$P(x_1 | c_1, b_2) = \alpha(0.4)(0.6)[(0.25)(0.7) + (025)(0.8) + (0.25)(0.1) + (0.25)(0.3)]1.0$$

$$P(x_1 | c_1, b_2) = \alpha 0.114$$

$$P(x_2 | c_1, b_2) = \alpha 0.066$$

classificação: salmon!

### Naive Bayes

- relações de dependência entre atributos desconhecidas
- neste caso assume-se independência condicional

$$P(\mathbf{x} \mid \mathbf{a}, \mathbf{b}) = P(\mathbf{x} \mid \mathbf{a}) P(\mathbf{x} \mid \mathbf{b})$$

## 9-Resumo

- Teoria Bayesiana de decisão é simples
- Regras de decisão
  - minimizar risco: ação que minimiza risco condicional
  - minimizar Pr[erro]: estado que maximiza densidade a posteriori  $P(\omega_i|x)$
- Superfícies decisão hiperquadráticas no caso Gaussiano
- Redes: relações dependência/independência entre variáveis

### Observação

Este material refere-se às notas de aula do curso CT 720 Tópicos Especiais em Aprendizagem de Máquina e Classificação de Padrões da Faculdade de Engenharia Elétrica e de Computação da Unicamp e do Centro Federal de Educação Tecnológica do Estado de Minas Gerais. Não substitui o livro texto, as referências recomendadas e nem as aulas expositivas. Este material não pode ser reproduzido sem autorização prévia dos autores. Quando autorizado, seu uso é exclusivo para atividades de ensino e pesquisa em instituições sem fins lucrativos.

ProfFernandoGomide ©DCA-FEEC-Unicamp