

# Curso 2 – CD, AM e DM



## **CLASSIFICADOR BAYESIANO** ***NAIVE-BAYES CLASSIFIER***

PROFA. ROSELI AP. FRANCELIN ROMERO

SCC – ICMC - USP



# Aprendizado Bayesiano



## CLASSIFICADORES BAYESIANO

**Aprendizado  
Supervisionado  
de  
Classificadores  
Bayesiano**

**Aprendizado  
Não Supervisionado  
de  
Classificadores  
Bayesiano**



# Classificação de Padrões

Suponha que você está para testemunhar um evento.  
O evento pertencerá à:

- classe  $\omega_1$  com probabilidade  $P(\omega_1)$
- classe  $\omega_2$  com probabilidade  $P(\omega_2)$
- classe  $\omega_n$  com probabilidade  $P(\omega_n)$

Suponha que você deve prever a classe:

- Você paga R\$ 1,00 se você estiver errado
- Você não paga nada se estiver certo.

Questões:

- Qual deve ser sua estratégia ótima?
- Qual será o seu custo esperado?



# Considerando dados observados



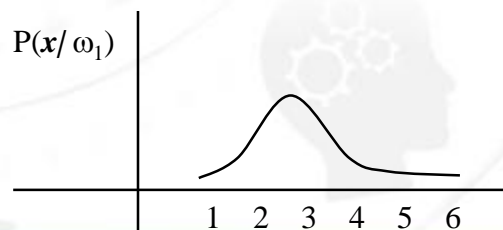
Suponha que se deseja construir um SISTEMA AUTOMÁTICO para apanhar **batatas**. Toda vez que um objeto toca o sensor debaixo do trator ele deve decidir se pertence à:

$\omega_1$  **batata** com probabilidade  $P(\omega_1)$

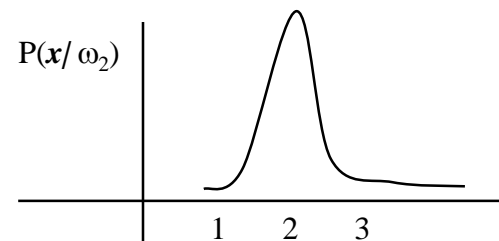
$\omega_2$  **pedra** com probabilidade  $P(\omega_2)$

$\omega_3$  **terrão** com probabilidade  $P(\omega_3)$

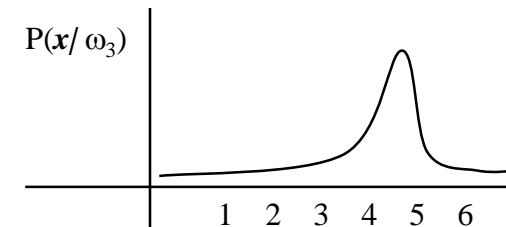
Suponha também que o sensor computa o diâmetro  $x$  do objeto e que o Instituto de Pesquisa da Batata forneceu as distribuições condicionais de  $x$  para cada classe.



BATATA



PEDRA



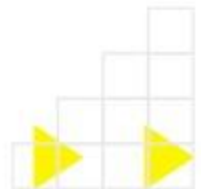
TERRÃO

# DECISÃO

- Conhece-se  $P(\omega_1)$ ,  $P(\omega_2)$ ,  $P(\omega_3)$  mais as distribuições  $P(\mathbf{x} / \omega_1)$ ,  $P(\mathbf{x} / \omega_2)$ ,  $P(\mathbf{x} / \omega_3)$ .
- Observa-se  $\mathbf{x}$ .
- Qual a classe de objetos escolhida?

## I - Máxima Probabilidade

- Escolher a classe  $\omega_i$  que maximiza  $P(\mathbf{x} / \omega_i)$ .
- Fácil de calcular.
- Qual é a objeção? (pode ocorrer erro! Porque se toma a probabilidade partindo-se de uma certa classe).



# DECISÃO

## II - Classificador Bayesiano Ótimo

O que devemos fazer para minimizar a chance de cometermos um erro?

- Escolher a classe  $\omega_i$  que tem a maior probabilidade dada  $\mathbf{x}$ .

$$\text{Escolha} = \arg_i \max P(\omega_i / \mathbf{x}).$$

Bayesiano Ótimo =

$$\arg_i \max P(\mathbf{x} / \omega_i) \cdot P(\omega_i)$$

Este é o Classificador Ótimo de Bayes.





# Batatas Multivariado

Suponha que temos 3 sensores

$$\left\{ \begin{array}{l} x_1 - \text{diâmetro} \\ x_2 - \text{altura} \\ x_3 - \text{massa} \end{array} \right.$$

e que temos um vetor  $\mathbf{x}$  observado

$$\text{Bayesiano Ótimo} = \arg_i \max_{\tilde{\omega}_i} P(\mathbf{x} / \omega_i) \cdot P(\omega_i)$$

Hipótese Comum:

Cada  $P(\mathbf{x} / \omega_i)$  segue **distribuição Gaussiana**.

Três Casos:

$P(\tilde{\mathbf{x}} / \omega_i)$  - Média  $\mu_i$ , variância  $\sigma^2$

$P(\tilde{\mathbf{x}} / \omega_i)$  - Média  $\mu_i$ , covariância  $\Sigma$ , arbitrária

$P(\tilde{\mathbf{x}} / \omega_i)$  - Média  $\mu_i$ , covariância  $\Sigma_i$ , diferente para classes diferentes

# Função Gaussiana

**Caso 1:** Todas componentes são independentes  $P(\mathbf{x}/\omega_i)$  tem média  $\mu_i$ . Cada componente de  $\mathbf{x}$  é independente de outras componentes e tem variância  $\sigma^2$ :

$$P(\tilde{\mathbf{x}} / \omega_i) = k \exp\left(-\frac{1}{2\sigma^2} \sum (\mathbf{x}_j - \mu_{ij})^2\right)$$

$$\begin{aligned} \text{Bayesiano Ótimo} &= \arg_i \max P(\tilde{\mathbf{x}} / \omega_i) \cdot P(\omega_i) = \\ &= \arg_i \max \left\{ k \exp\left(-\frac{1}{2\sigma^2} \sum (\mathbf{x}_j - \mu_{ij})^2\right) \cdot P(\omega_i) \right\} = \end{aligned}$$

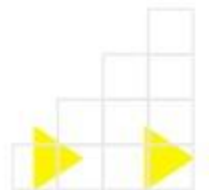


# Classificação de Padrões

- Suponha agora que voce nao conhece  
 $P(w_1) \ P(w_2) \ ... \ P(w_N) \ , \ \mu_1, \mu_2 \ ... \ \mu_N$

Mas, voce deseja estimar estes parametros dos dados.

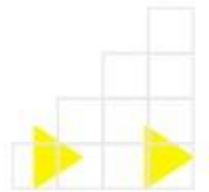
$\tilde{x}_1^{(1)} \ \tilde{x}_2^{(1)} \ ... \ \tilde{x}_N^{(1)}$	<b>Classe <math>w_1</math></b>
$\tilde{x}_1^{(2)} \ \tilde{x}_2^{(2)} \ ... \ \tilde{x}_N^{(2)}$	<b>Classe <math>w_2</math></b>
$\ddots$	
$\tilde{x}_1^{(M)} \ \tilde{x}_2^{(M)} \ ... \ \tilde{x}_N^{(M)}$	<b>Classe <math>w_N</math></b>





# Classificacao de Padroes

- Estimar  $P(w_i) = \frac{\text{numero de dados da classe } w_i}{\text{numero total de dados}}$
- Estima a media  $\mu_i$  = media de todos os pontos da classe  $w_i$





# Métodos de Aprendizado Bayesiano

- Calculam explicitamente probabilidades para hipóteses (Naïve Bayes Classificador).

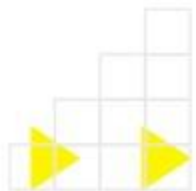
Mitchie et al. (1994) comparou o classificador Naïve Bayes com RN e DT.

- Eles fornecem uma perspectiva útil para compreensão dos algoritmos de aprendizado que não explicitamente manipulam probabilidades.



# Características dos Métodos de Aprendizado Bayesiano

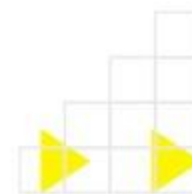
- Cada exemplo observado pode incrementalmente diminuir ou aumentar a probabilidade estimada que uma hipótese está correta.



# Características dos Métodos de Aprendizado Bayesiano



- Métodos Bayesiano podem acomodar hipóteses que contém previsões probabilísticas, tais como:  
“este paciente, com pneumonia, tem 93% de chance de cura”.
- Novas instâncias podem ser classificadas combinando as previsões de múltiplas hipóteses, ponderadas por “***suas probabilidades***”.
- Em métodos computacionais igualmente intratáveis, eles podem fornecer um padrão de tomada de decisão ótima.



# Características dos Métodos de Aprendizado Bayesiano

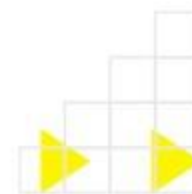


- Dificuldade 1:

Requerem o conhecimento de muitas probabilidades. Quando estas probabilidades não são conhecidas “a priori” elas são estimadas baseadas no: **conhecimento do problema, dados previamente disponíveis e hipóteses sobre a forma da distribuição fundamental dos dados.**

- Dificuldade 2:

Custo computacional requerido, mas pode ser reduzido significativamente.







# TEOREMA DE BAYES

Em **problemas de AM** estamos interessados em  $P(h | D)$ : probabilidade posteriori, probabilidade vale  $h$  dado o conjunto de treinamento observado  $D$ .

Teorema de Bayes:

$$P(h | D) = \frac{P(D|h) P(h)}{P(D)}$$

Em muitos casos o aprendiz considera algum conjunto de hipóteses candidatas  $H$  e está interessado em encontrar a hipótese mais provável  $h \in H$  dado o conjunto de dados observado  $D$  ( ou no mínimo a hipótese mais provável, se existirem várias).



# TEOREMA DE BAYES

Tal hipótese é chamada uma Maximum A Posteriori (MAP) hipótese.

$$\begin{aligned} h_{\text{MAP}} &= \arg_{h \in H} \max P(h | D) = \\ &= \arg_{h \in H} \max \frac{P(D | h) P(h)}{P(D)} \longrightarrow \text{É independente de } h \\ &= \arg_{h \in H} \max P(D | h) P(h) \end{aligned}$$

Em alguns casos, assumiremos que toda hipótese em  $H$  é igualmente provável, isto é:

$P(h_i) = P(h_j)$  para todos  $h_i$  e  $h_j$  em  $H$  então a equação anterior fica:



# TEOREMA DE BAYES

$$h_{ML} = \arg_{h \in H} \max P(D | h)$$

Maximum likelihood (Probabilidade Maxima)

## No enfoque de ML

**D** - exemplos de treinamento de alguma função alvo.

**H** - como o espaço das funções alvo candidatas.

# EXEMPLO

Paciente tem câncer ou não?

Um paciente faz um teste de laboratório e o resultado volta positivo. O teste devolve um resultado positivo correto em só 98% dos casos nos quais a doença está realmente presente, e um resultado negativo correto em 97% dos casos nos quais a doença não está presente. Além disso, 0.008 da população inteira tem este câncer.

$$P(\text{câncer}) = 0.008$$

$$P(\neg \text{câncer}) = 0.992$$

$$P(+ | \text{câncer}) = 0.98$$

$$P(- | \text{câncer}) = 0.02$$

$$P(+ | \neg \text{câncer}) = 0.03$$

$$P(- | \neg \text{câncer}) = 0.97$$

$$P(+ | \text{câncer}) \cdot P(\text{câncer}) = (0.98) \cdot (0.008) = 0.0078$$

$$P(+ | \neg \text{câncer}) \cdot P(\neg \text{câncer}) = (0.03) \cdot (0.992) = 0.0298$$

$$h_{\text{MAP}} = \neg \text{câncer}$$



# Classificador Bayesiano Naive

Está entre um dos melhores classificadores (árvores de decisão, NN, KNN)

Quando usar:

- Conjunto de treinamento grande.
- Atributos são condicionalmente independentes.

Aplicações bem sucedidas:

- Diagnósticos
- Classificação de textos em documentos





# Classificador Bayesiano Naive

Seja:

$$f: X \rightarrow V$$

$$x = \langle a_1, a_2, \dots, a_n \rangle$$

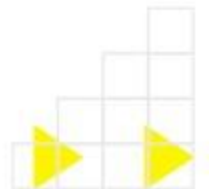
Qual é o mais provável valor de  $f(x)$  ?

$$v_{MAP} = \arg_{v_j \in V} \max P(v_j | a_1, a_2, \dots, a_n)$$

$$v_{MAP} = \arg_{v_j \in V} \max \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

$$v_{MAP} = \arg_{v_j \in V} \max P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$

Hipótese de Naïve Bayes:  $P(a_1, a_2, \dots, a_n | v_j) = \prod_i P(a_i | v_j)$





# Classificador Bayesiano Naive

Classificador Bayesiano Naïve:

$$V_{NB} = \arg_{v_j \in V} \max P(v_j) \prod_i P(a_i | v_j)$$

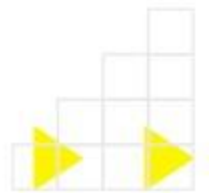
## EXEMPLO:

Considere o exemplo “Play Tennis” e a instância:

**<Outlook = sol,Temp=fria,Um=alta,vento=forte>**

Queremos:

$$V_{NB} = \arg_{v_j \in V} \max P(v_j) \prod_i P(a_i | v_j) =$$





# Classificador Bayesiano Naive

$$\Rightarrow P(\text{sim}) P(\text{sol} | \text{sim}) P(\text{frio} | \text{sim}) P(\text{alta} | \text{sim}) P(\text{forte} | \text{sim}) = 0.0053$$

$$\Rightarrow P(\text{n\~ao}) P(\text{sol} | \text{n\~ao}) P(\text{frio} | \text{n\~ao}) P(\text{alta} | \text{n\~ao}) P(\text{forte} | \text{n\~ao}) = 0.0206$$

$$\rightarrow V_{NB} = n$$

**OBS: Cap.6** - T. Mitchell para ver aplicação de busca de texto em documentos da Web.



# REFERÊNCIAS

- MITCHELL, Tom M. Machine Learning: McGraw-Hill, 1997.
- WEKA, Weka 3: Data Mining Software in Java, Disponível em [www.cs.waikato.ac.nz/ml/weka/](http://www.cs.waikato.ac.nz/ml/weka/), Acesso em: mar. 2010.

