#### APRENDIZADO BAYESIANO

APRENDIZADO DE MÁQUINA

PROFA. Roseli Ap. Francelin Romero

# Fórmulas Básicas para Probabilidades

Regra Produto: probabilidade  $P(A \land B)$  de uma conjunção de dois eventos  $A \in B$ :

$$P(A \land B) = P(A \mid B) P(B) = P(B \mid A) P(A)$$

Regra Soma: probabilidade  $P(A \lor B)$  de uma união de dois eventos  $A \in B$ :

$$P(A \lor B) = P(A) + P(B) - P(A \land B)$$

Teorema da probabilidade total: se eventos  $A_1, \ldots, A_n$  são mutualmente exclusivos com  $\sum_{i=1}^{n} P(A_i) = 1$ , então:

$$P(B) = \sum_{i=1}^{n} P(B \mid A_i) P(A_i)$$

# 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:

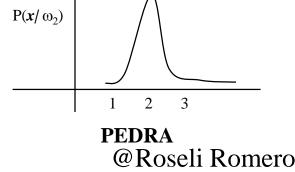
- $\omega_1$  batata com probabilidade  $P(\omega_1)$
- $\omega_2$  *pedra* com probabilidade  $P(\omega_2)$
- $\omega_3$  terrão com probabilidade  $P(\omega_3)$

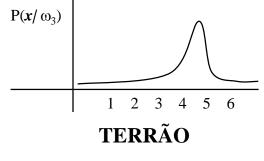
Suponha também que o sensor computa o diâmetro x do objeto e que o Instituto de Pesquisa da Batata forneceu as distribuições condicionais de x para cada classe.

P(x/ω<sub>1</sub>)

1 2 3 4 5 6

BATATA





### **DECISÃO**

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

#### I - Máxima Probabilidade

- Escolher a classe  $\omega_i$  que maximiza  $P(x/\omega_i)$ .
- 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 x.

Escolha = 
$$arg_i max P(\omega_i/x)$$
.  
Bayesiano Ótimo =  $arg_i max P(x/\omega_i).P(\omega_i)$ 

Este é o Classificador Ótimo de Bayes.

#### **Batatas Multivariado**

Suponha que temos 3 sensores 
$$\begin{cases} x_1 - \text{diâmetro} \\ x_2 - \text{altura} \\ x_3 - \text{massa} \end{cases}$$

e que temos um vetor x observado

Bayesiano Ótimo =  $\arg_{i} \max P(\tilde{x} / \omega_{i})$ .  $P(\omega_{i})$ 

Hipótese Comum:

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

Três Casos:

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

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

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

## Função Gaussiana

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

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

Bayesiano Ótimo = 
$$\arg_{i}\max P(\boldsymbol{x}/\omega_{i})$$
.  $P(\omega_{i}) = \arg_{i}\max\{k \exp(-\frac{1}{2\sigma^{2}}) (\boldsymbol{x}_{j} - \mu_{i})^{2})$ .  $P(\omega_{i}) = \frac{1}{2\sigma^{2}}$ 

#### Caso 1

= 
$$\underset{2\sigma^2}{\text{arg}_i \text{max}} \frac{1}{2\sigma^2} \sum_{j} (x_j - \mu_{ij})^2 + \log P(\omega_i) =$$

$$= \underset{j}{\operatorname{arg_imin}} \quad \sum_{j} (x_j - \mu_{ij})^2 - 2\sigma^2 \log P(\omega_i) = \frac{2\sigma^2}{2\sigma^2}$$

= arg<sub>i</sub>min 
$$\sum_{i} (x_j - \mu_{ij})^2 - 2\sigma^2 \log P(\omega_i)$$

#### Caso 1

- Caso duas classes
- =  $\underset{\sim}{\text{arg}_{i}} \min \left( (x \mu_{i})^{2} 2\sigma^{2} \log P(\omega_{i}) \right) =$
- = arg<sub>i</sub>min  $(x_i x_i 2x_i \mu_i + \mu_i \mu_i 2\sigma^2 \log P(\omega_i)) =$
- =  $arg_i min \left( -2 x \mu_i + c_i \right)$

Se - 
$$2 \times \mu_1 + c_1 < -2 \times \mu_2 + c_2 \rightarrow Escolha \omega_1 \Leftrightarrow$$

#### Caso 1

- $\Leftrightarrow$  Se  $c_1 c_2 < 2 (\mu_1 \mu_2) x \rightarrow Escolha <math>\omega_1$
- ⇔ A regra de decisão é:
- "Se  $\omega x > threshold$ " onde  $\omega = 2 (\mu_1 \mu_2)$ e  $threshold = c_1 - c_2$
- Portanto a decisão ótima é de um CLASSIFICADOR LINEAR! Perceptrons são corretos!
- OBS.: A regra do Perceptron pode ser obtida do classificador ótimo de Bayes.

# Caso 2 - Hipótese mais fraca

Agora,  $P(x/\omega_i)$  gaussiana, média  $\mu_i$  e covariância arbitrária  $\Sigma$ . Temos que a mesma regra ocorre, mas numa medida de distancia diferente:

Dist 
$$(\underline{x}, \underline{\mu}_i) = (\underline{x} - \underline{\mu}_i)^T \sum_{i=1}^{T-1} (\underline{x} - \underline{\mu}_i)$$

Se todos os  $P(\omega_i)_S$  são iguais  $\Rightarrow$  método do vizinho mais próximo (KNN)

Ainda usa regiões de decisão linear.

# Caso 3 - Hipótese ainda mais fraca

 $P(x/\omega_i)$  gaussiana, média  $\mu_i$  e covariância  $\Sigma_i \rightarrow$  para diferentes classes a variância pode ser diferente.

Ainda é fácil calcular a decisão ótima

$$\underset{\sim}{\operatorname{arg_{i}max}}(P(\omega_{i}/x))$$

mas as regiões de decisão não são mais lineares.

#### Classificação de Padroes

Suponha agora que voce nao conhece

$$P(w_1) P(w_2) \dots P(w_N) , \mu_1, \mu_2 \dots \mu_N$$

Mas, voce deseja estimar estes parametros dos dados.

$$x_1^{(1)} x_2^{(1)} ... x_1^{(1)} N$$
 Classe  $w_1$   $x_1^{(2)} x_2^{(2)} ... x_1^{(2)} N$  Classe  $w_2$ 

 $\mathbf{x_1}^{(\mathbf{M})} \ \mathbf{x_2}^{(\mathbf{M})} \ \dots \ \mathbf{x^{(\mathbf{M})}_N}$  Classe  $\mathbf{w_N}$ 

#### Classificação de Padroes

Estimar  $P(w_i) = \underline{\text{numero de dados da classe } w_i}$ 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.
- Conhecimento "priori" pode ser combinado com o dado observado para determinar a probabilidade final de uma hipótese. Em Aprendizado Bayesiano, conhecimento a priori, pode ser fornecido:
  - Dando uma probabilidade "a priori" para cada hipótese candidata.
  - Distribuição de probabilidade sobre os dados para cada hipótese possível.

# 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 pode ser reduzido significantemente.

#### TEOREMA DE BAYES

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

Teorema de Bayes:  $P(h|D) = \frac{P(D|n) P(n)}{P(D)}$ 

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

#### TEOREMA DE BAYES

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

$$\begin{aligned} h_{MAP} &= arg_{h \in H} max \ P(h|D) = \\ &= arg_{h \in H} max \ \underline{P(D|h) \ P(h)} = \underbrace{ \ \acute{E} \ independente \ de \ h} \\ &\underline{P(D)} \end{aligned}$$

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

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.

$$P(cancer) = 0.008 \qquad P(\neg cancer) = 0.992$$

$$P(+|cancer| = 0.98)$$
  $P(-|cancer| = 0.02)$ 

$$P(+|\neg cancer) = 0.03$$
  $P(-|\neg cancer) = 0.97$ 

$$P(+|cancer) \cdot P(cancer) = (0.98) \cdot (0.008) = 0.0078$$

$$P(+|\neg cancer) \cdot P(\neg cancer) = (0.03) \cdot (0.992) = 0.0298$$

$$h_{MAP} = \neg c \hat{a} n c e r$$

# Classificação mais Provável de Novas Instâncias

- Até agora nós buscamos a mais provável hipótese dado o conjunto  $\bf D$  (i.e.  $\bf h_{MAP}$ )
- Dado nova instância **x**, qual é a sua classificação mais provável?
  - $\mathbf{h}_{\mathbf{MAP}}(\mathbf{x})$  não é a classificação mais provável.

#### Considere por exemplo:

- três hipóteses:  $P(h_1|D)=0.4, P(h_2|D)=0.3, P(h_3|D)=0.3$
- Dado a nova instância x:  $h_1(x) = +, h_2(x) = -, h_3(x) = -$
- Qual é a mais provável classificação de x?
  - $\mathbf{p}_{+}(\mathbf{x}) = \mathbf{0.4}$ ,  $\mathbf{p}_{-}(\mathbf{x}) = \mathbf{0.6}$ , portanto é mais provável que  $\mathbf{x}$  seja -

Neste caso, é diferente da classificação gerada pela  $\mathbf{h}_{\mathbf{MAP}}$ 

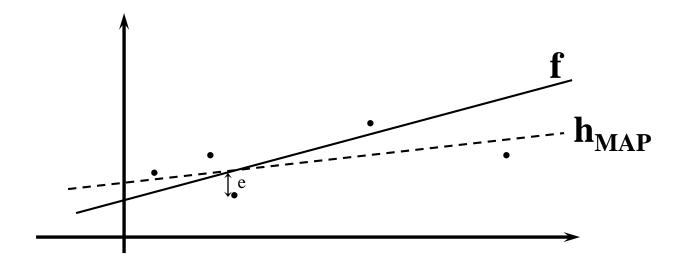
# Classificador Bayesiano Ótimo

$$\arg_{\mathbf{v_j} \in \mathbf{V}} \max \sum_{\mathbf{h_i} \in \mathbf{H}} \mathbf{P}(\mathbf{v_j} / \mathbf{h_i}) \cdot \mathbf{P}(\mathbf{h_i} / D)$$

#### **EXEMPLO:**

$$\begin{aligned} & P(h_1/D) = 0.4, & P(-/h_1) = 0, & P(+/h_1) = 1, \\ & P(h_2/D) = 0.3, & P(-/h_2) = 1, & P(+/h_2) = 0, \\ & P(h_3/D) = 0.3, & P(-/h_3) = 1, & P(+/h_3) = 0, \\ & P(-/h_3) = 1, & P(+/h_3) = 0, \\ & P(-/h_3) = 1, & P(-/h_3) = 0, \\ & P(-/h_3) = 1, & P(-/h_3) = 0, \\ & P(-/h_3) = 0, & P(-/h_3) = 0.4 \\ & P(-/h_3) = 0, & P(-/h_3) = 0.4 \\ & P(-/h_3) = 0, & P(-/h_3) = 0.4 \\ & P(-/h$$

# Aprendizado de uma Função Real



Considere exemplos de treinamento  $\langle x_i, d_i \rangle$ , onde  $d_i$  é o ruído dado por:

$$\mathbf{d_i} = \mathbf{f}(\mathbf{x_i}) + \mathbf{e_i}$$

onde **e**<sub>i</sub> é uma variável aleatória, independente para @Roseli Romero

# Aprendizado de uma Função Real

Cada  $x_i$  de acordo com alguma distribuição Gaussiana com média = 0. Então,

$$\mathbf{h}_{\mathrm{ML}} = \arg_{\mathbf{h} \in \mathbf{H}} \min \sum_{i=1}^{m} (\mathbf{d}_{i} - \mathbf{h}(\mathbf{x}_{i}))^{2}$$

#### Demonstração:

$$h_{ML} = arg_{h \in H} max \ p(D|h) = arg_{h \in H} max \ \underset{i=1}{\overset{m}{\pi}} p(d_i|h) =$$

$$= arg_{h \in H} max \int_{i=1}^{m} \frac{1}{\sqrt{2\sigma^2}} e^{-1/2((d_i - h(x_i))|\sigma)^2} =$$

Maximizando o logaritmo natural:

$$\mathbf{h}_{\mathrm{ML}} = \mathbf{arg}_{\mathbf{h} \in \mathbf{H}} \mathbf{max} \sum_{i=1}^{-1/2} ((\mathbf{d}_{i} - \mathbf{h}(\mathbf{x}_{i}))/\sigma)^{2} =$$
@Roseli Romero

# Aprendizado de uma Função Real

$$= \arg_{h \in H} \max \sum_{i=1}^{m} -(d_i - h(x_i))^2 =$$

= 
$$\underset{i=1}{\operatorname{arg}} \min \sum_{i=1}^{m} (\mathbf{d_i} - \mathbf{h}(\mathbf{x_i}))^2$$

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

Seja: 
$$\begin{aligned} \mathbf{f} \colon \mathbf{X} &\to \mathbf{V} \\ \mathbf{x} &= <\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n > \end{aligned}$$
 Qual é o mais provável valor de  $\mathbf{f}(\mathbf{x})$ ? 
$$\mathbf{v}_{\text{MAP}} &= \mathbf{arg}_{\mathbf{v}_j \in \mathbf{V}} \mathbf{max} \; \mathbf{P}(\mathbf{v}_j \, | \, \mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, ) \\ \mathbf{v}_{\text{MAP}} &= \mathbf{arg}_{\mathbf{v}_j \in \mathbf{V}} \mathbf{max} \; \underline{\mathbf{P}}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{v}_j) \\ &\qquad \qquad \qquad \underline{\mathbf{P}}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{a}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{a}_i \, | \, \mathbf{v}_j) \\ \mathbf{P}(\mathbf{a}_1, \, \mathbf{a}_2, \, \ldots, \, \mathbf{v}_n \, | \, \mathbf{v}_j \, ) \; \mathbf{P}(\mathbf{v}_j \, | \, \mathbf{v}_j)$$

Classificador Bayesiano Naïve:

$$\mathbf{V}_{NB} = \arg_{\mathbf{v_j} \in \mathbf{V}} \max \mathbf{P}(\mathbf{v_j}) \prod_{i} \mathbf{P}(\mathbf{a_i} | \mathbf{v_j})$$

**EXEMPLO:** 

Considere o exemplo "Play Tennis" e a instância:

<Outlook = sunny,Temp=cool,Hum=high,wind=strong>
Queremos:

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

- ⇒P(yes) P(sunny|yes) P(cool|yes) P(high|yes) P(strong|yes)= = 0.0053
- $\Rightarrow$ P(no) P(sunny|no) P(cool|no) P(high|no) P(strong|no)= = 0.0206

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

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