

APRENDIZADO BAYESIANO

SCC0276 – APRENDIZADO DE MÁQUINA

PROFA. Roseli Ap. Francelin Romero

Fórmulas Básicas para Probabilidades

- Regra Produto: probabilidade $P(\mathbf{A} \wedge \mathbf{B})$ de uma conjunção de dois eventos \mathbf{A} e \mathbf{B} :

$$P(\mathbf{A} \wedge \mathbf{B}) = P(\mathbf{A} \mid \mathbf{B}) P(\mathbf{B}) = P(\mathbf{B} \mid \mathbf{A}) P(\mathbf{A})$$

- Regra Soma: probabilidade $P(\mathbf{A} \vee \mathbf{B})$ de uma união de dois eventos \mathbf{A} e \mathbf{B} :

$$P(\mathbf{A} \vee \mathbf{B}) = P(\mathbf{A}) + P(\mathbf{B}) - P(\mathbf{A} \wedge \mathbf{B})$$

- Teorema da probabilidade total: se eventos $\mathbf{A}_1, \dots, \mathbf{A}_n$ são mutualmente exclusivos com $\sum_{i=1}^n P(\mathbf{A}_i) = 1$, então:

$$P(\mathbf{B}) = \sum_{i=1}^n P(\mathbf{B} \mid \mathbf{A}_i) P(\mathbf{A}_i)$$

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 ω_1 com probabilidade $P(\omega_1)$
 - classe ω_2 com probabilidade $P(\omega_2)$
 - classe ω_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

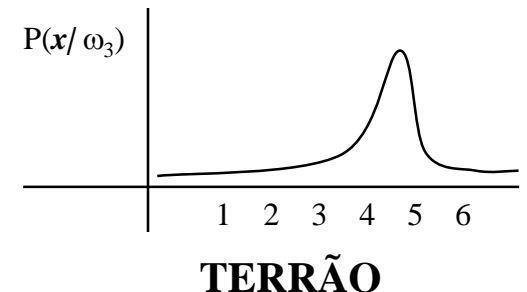
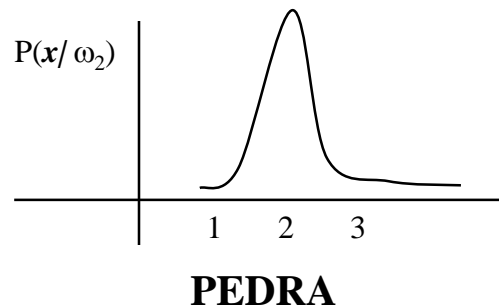
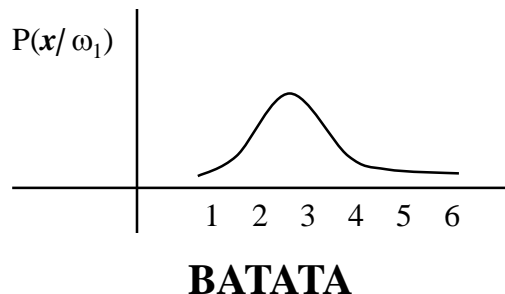
Suponha que se deseja construir um SISTEMA AUTOMÁTICO para apanhar *batatas*. Toda vez que um objeto toca o sensor debaixo do trator ele deve decidir se pertence à:

ω_1 *batata* com probabilidade $P(\omega_1)$

ω_2 *pedra* com probabilidade $P(\omega_2)$

ω_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.



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 ω_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 ω_i que tem a maior probabilidade dada x .

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

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

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

$$\text{Bayesiano Ótimo} = \arg_i \max P(\mathbf{x} / \omega_i) \cdot P(\omega_i)$$

Hipótese Comum:

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

Três Casos:

$P(\mathbf{x} / \omega_i)$ - Média μ_i , variância σ^2

$P(\mathbf{x} / \omega_i)$ - Média μ_i , covariância Σ , arbitrária

$P(\mathbf{x} / \omega_i)$ - Média μ_i , covariância Σ_i , diferente para classes diferentes

Caso 1: Todas componentes são independentes

$P(\mathbf{x}/\omega_i)$ tem média μ_i . Cada componente de \mathbf{x} é independente de outras componentes e tem variância σ^2

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

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

$$= \arg_i \max \frac{-1}{2\sigma^2} \sum_j (\mathbf{x}_j - \mu_{ij})^2 + \log P(\omega_i) =$$

$$= \arg_i \min \frac{\sum_j (\mathbf{x}_j - \mu_{ij})^2 - 2\sigma^2 \log P(\omega_i)}{2\sigma^2} =$$

$$= \arg_i \min \sum_j (\mathbf{x}_j - \mu_{ij})^2 - 2\sigma^2 \log P(\omega_i)$$

■ Caso duas classes

$$= \arg_i \min ((\tilde{x} - \tilde{\mu}_i)^2 - 2\sigma^2 \log P(\omega_i)) =$$

$$= \arg_i \min (\tilde{x} \tilde{x} - 2 \tilde{x} \tilde{\mu}_i + \tilde{\mu}_i \tilde{\mu}_i - 2\sigma^2 \log P(\omega_i)) =$$

$$= \arg_i \min (- 2 \tilde{x} \tilde{\mu}_i + c_i)$$

$$\text{Se } - 2 \tilde{x} \tilde{\mu}_1 + c_1 < - 2 \tilde{x} \tilde{\mu}_2 + c_2 \rightarrow \text{Escolha } \omega_1 \Leftrightarrow$$

\Leftrightarrow Se $c_1 - c_2 < 2 (\mu_1 - \mu_2) x \rightarrow$ Escolha ω_1

\Leftrightarrow 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.

Hipótese mais fraca

Agora, $P(\underline{x}/\omega_i)$ gaussiana, média $\underline{\mu}_i$ e covariância arbitrária Σ . Temos que a mesma regra ocorre, mas numa medida de distancia diferente:

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

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

Ainda usa regiões de decisão linear.

Hipótese ainda mais fraca

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

Ainda é fácil calcular a decisão ótima

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

mas as regiões de decisão não são mais lineares.

Classificacao 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.

$$\mathbf{x}_1^{(1)} \ \mathbf{x}_2^{(1)} \ \dots \ \mathbf{x}_N^{(1)}$$

Classe w_1

$$\mathbf{x}_1^{(2)} \ \mathbf{x}_2^{(2)} \ \dots \ \mathbf{x}_N^{(2)}$$

Classe w_2

..

$$\mathbf{x}_1^{(M)} \ \mathbf{x}_2^{(M)} \ \dots \ \mathbf{x}_N^{(M)}$$

Classe w_N

Classificacao de Padroes

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

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 prior, 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 significativamente.

TEOREMA DE BAYES

Em problemas de ML estamos interessados em $\mathbf{P}(\mathbf{h}|\mathbf{D})$:
probabilidade a posteriori, probabilidade vale \mathbf{h} dado o
conjunto de treinamento observado \mathbf{D} .

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

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.

$$\begin{aligned} \mathbf{h}_{\text{MAP}} &= \arg_{\mathbf{h} \in \mathbf{H}} \max \mathbf{P}(\mathbf{h}|\mathbf{D}) = \\ &= \arg_{\mathbf{h} \in \mathbf{H}} \max \frac{\mathbf{P}(\mathbf{D}|\mathbf{h}) \mathbf{P}(\mathbf{h})}{\mathbf{P}(\mathbf{D})} = \end{aligned}$$

É independente de \mathbf{h}

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

Em alguns casos, assumiremos que toda hipótese em \mathbf{H} é igualmente provável, isto é:

$\mathbf{P}(\mathbf{h}_i) = \mathbf{P}(\mathbf{h}_j)$ para todos \mathbf{h}_i e \mathbf{h}_j em \mathbf{H} então a equação anterior fica:

TEOREMA DE BAYES

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

 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.

$$P(\text{câncer}) = 0.008$$

$$P(\neg \text{câncer}) = 0.992$$

$$P(+|\text{câncer}) = 0.98$$

$$P(-|\text{câncer}) = 0.02$$

$$P(+|\neg \text{câncer}) = 0.03$$

$$P(-|\neg \text{câncer}) = 0.97$$

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

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

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

Classificação mais Provável de Novas Instâncias

Até agora nós buscamos a mais provável hipótese dado o conjunto \mathbf{D} (i.e. \mathbf{h}_{MAP})

Dado nova instância \mathbf{x} , qual é a sua classificação mais provável?

$\mathbf{h}_{\text{MAP}}(\mathbf{x})$ não é a classificação mais provável.

Considere por exemplo:

- ♦ três hipóteses: $\mathbf{P}(\mathbf{h}_1|\mathbf{D})=0.4$, $\mathbf{P}(\mathbf{h}_2|\mathbf{D})=0.3$, $\mathbf{P}(\mathbf{h}_3|\mathbf{D})=0.3$
- ♦ Dado a nova instância \mathbf{x} : $\mathbf{h}_1(\mathbf{x})=+$, $\mathbf{h}_2(\mathbf{x})=-$, $\mathbf{h}_3(\mathbf{x})=-$
- ♦ Qual é a mais provável classificação de \mathbf{x} ?

$\mathbf{p}_+(\mathbf{x})=0.4$, $\mathbf{p}_-(\mathbf{x})=0.6$, portanto é mais provável que \mathbf{x} seja -

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

Classificador Bayesiano Ótimo

$$\arg_{\mathbf{v}_j \in V} \max \sum_{\mathbf{h}_i \in H} P(\mathbf{v}_j / \mathbf{h}_i) \cdot P(\mathbf{h}_i / D)$$

EXEMPLO:

$$P(\mathbf{h}_1 / D) = 0.4, \quad P(- / \mathbf{h}_1) = 0, \quad P(+ / \mathbf{h}_1) = 1,$$

$$P(\mathbf{h}_2 / D) = 0.3, \quad P(- / \mathbf{h}_2) = 1, \quad P(+ / \mathbf{h}_2) = 0,$$

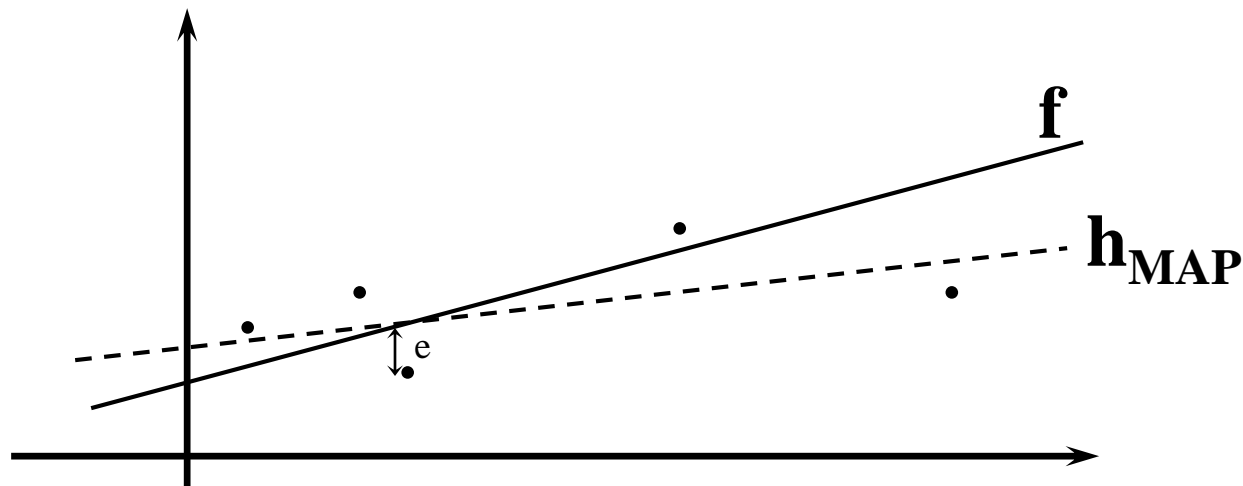
$$P(\mathbf{h}_3 / D) = 0.3, \quad P(- / \mathbf{h}_3) = 1, \quad P(+ / \mathbf{h}_3) = 0,$$

Portanto, $\sum_{\mathbf{h}_i \in H} P(+ / \mathbf{h}_i) \cdot P(\mathbf{h}_i / D) = 0.4$

$$\sum_{\mathbf{h}_i \in H} P(- / \mathbf{h}_i) \cdot P(\mathbf{h}_i / D) = 0.6$$

Portanto, $\arg_{\mathbf{v}_j \in V} \max \sum_{\mathbf{h}_i \in H} P(\mathbf{v}_j / \mathbf{h}_i) \cdot P(\mathbf{h}_i / D) = -$

Aprendizado de uma Função Real



Considere exemplos de treinamento $\langle \mathbf{x}_i, \mathbf{d}_i \rangle$, onde \mathbf{d}_i é o ruído dado por:

$$\mathbf{d}_i = \mathbf{f}(\mathbf{x}_i) + \mathbf{e}_i$$

onde \mathbf{e}_i é uma variável aleatória, independente para

Aprendizado de uma Função Real

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

$$\mathbf{h}_{\text{ML}} = \arg_{\mathbf{h} \in H} \min \sum_{i=1}^m (d_i - h(x_i))^2$$

Demonstração:

$$\begin{aligned} \mathbf{h}_{\text{ML}} &= \arg_{\mathbf{h} \in H} \max p(\mathbf{D}|\mathbf{h}) = \arg_{\mathbf{h} \in H} \max \prod_{i=1}^m p(d_i | \mathbf{h}) = \\ &= \arg_{\mathbf{h} \in H} \max \prod_{i=1}^m \frac{1}{\sqrt{2\sigma^2}} e^{-1/2((d_i - h(x_i))/\sigma)^2} = \end{aligned}$$

Maximizando o logaritmo natural:

$$\mathbf{h}_{\text{ML}} = \arg_{\mathbf{h} \in H} \max \sum_{i=1}^m -1/2((d_i - h(x_i))/\sigma)^2 =$$

Aprendizado de uma Função Real

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

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

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: $\mathbf{f}: \mathbf{X} \rightarrow \mathbf{V}$

$$\mathbf{x} = \langle \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \rangle$$

Qual é o mais provável valor de $\mathbf{f}(\mathbf{x})$?

$$\mathbf{v}_{\text{MAP}} = \mathbf{arg}_{\mathbf{v}_j \in \mathbf{V}} \max \mathbf{P}(\mathbf{v}_j | \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)$$

$$\mathbf{v}_{\text{MAP}} = \mathbf{arg}_{\mathbf{v}_j \in \mathbf{V}} \max \frac{\mathbf{P}(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n | \mathbf{v}_j) \mathbf{P}(\mathbf{v}_j)}{\mathbf{P}(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)}$$

$$\mathbf{v}_{\text{MAP}} = \mathbf{arg}_{\mathbf{v}_j \in \mathbf{V}} \max \mathbf{P}(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n | \mathbf{v}_j) \mathbf{P}(\mathbf{v}_j)$$

Hipotese de Naïve Bayes: $\mathbf{P}(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n | \mathbf{v}_j) = \prod_i \mathbf{P}(\mathbf{a}_i | \mathbf{v}_j)$

Classificador Bayesiano Naive

Classificador Bayesiano Naïve:

$$\mathbf{V}_{\text{NB}} = \arg_{\mathbf{v}_j \in \mathbf{V}} \max \mathbf{P}(\mathbf{v}_j) \prod_i \mathbf{P}(\mathbf{a}_i | \mathbf{v}_j)$$

EXEMPLO:

Considere o exemplo “Play Tennis” e a instância:

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

Queremos:

$$\mathbf{V}_{\text{NB}} = \arg_{\mathbf{v}_j \in \mathbf{V}} \max \mathbf{P}(\mathbf{v}_j) \prod_i \mathbf{P}(\mathbf{a}_i | \mathbf{v}_j) =$$

Classificador Bayesiano Naive

$$\Rightarrow P(\text{yes}) P(\text{sunny}|\text{yes}) P(\text{cool}|\text{yes}) P(\text{high}|\text{yes}) P(\text{strong}|\text{yes}) = 0.0053$$

$$\Rightarrow P(\text{no}) P(\text{sunny}|\text{no}) P(\text{cool}|\text{no}) P(\text{high}|\text{no}) P(\text{strong}|\text{no}) = 0.0206$$

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

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