

Aprendizagem de Máquina

Alessandro L. Koerich

Programa de Pós-Graduação em Engenharia Elétrica
Universidade Federal do Paraná (UFPR)

Aprendizagem Bayesiana

Plano de Aula

- Introdução
- Teorema de Bayes
- Classificador Ótimo de Bayes
- Classificador Naïve Bayes
- Exemplos
- Resumo



Referências

- Duda R., Hart P., Stork D. Pattern Classification 2ed. Willey Interscience, 2002. Capítulos 2 & 3
- Mitchell T. Machine Learning. WCB McGraw–Hill, 1997. Capítulo 6.
- Theodoridis S., Koutroumbas K. Pattern Recognition. Academic Press, 1999. Capítulo 2

Introdução

- O pensamento Bayesiano fornece uma abordagem probabilística para aprendizagem
- Está baseado na suposição de que as quantidades de interesse são reguladas por **distribuições de probabilidade.**
- Distribuição de probabilidade: é uma função que descreve a probabilidade de uma variável aleatória assumir certos valores.

Introdução

- Decisões ótimas podem ser tomadas com base nestas probabilidades conjuntamente com os dados observados.
- Fornece a base para algoritmos de aprendizagem que manipulam probabilidades, bem como para outros algoritmos que não manipulam probabilidades explicitamente.

Introdução

- Os métodos Bayesianos são importantes por dois motivos:
 1. Fornecem algoritmos práticos de aprendizagem:
 - *Naïve Bayes*
 - Redes Bayesianas
 - Combinam conhecimento a priori com os dados observados
 - Requerem probabilidades a priori
 2. Fornecem uma estrutura conceitual útil:
 - “Norma de Ouro” para avaliar outros algoritmos de aprendizagem. Norma de Ouro → menor erro possível

Características da Aprendizagem Bayesiana

- Cada exemplo de treinamento pode decrementar ou incrementar a probabilidade de uma hipótese ser correta.
- Conhecimento *a priori* pode ser combinado com os dados observados para determinar a probabilidade de uma hipótese.
- Métodos Bayesianos podem acomodar hipóteses que fazem previsões probabilísticas. Ex.: o paciente tem uma chance de 93% de possuir a doença.
- Novas instâncias podem ser classificadas combinando a probabilidade de múltiplas hipóteses ponderadas pelas suas probabilidades.

Dificuldades Práticas

- Métodos Bayesianos requerem o conhecimento inicial de várias probabilidades.
 - Quando não conhecidas, podem ser estimadas:
 - a partir de conhecimento prévio
 - dados previamente disponíveis
 - suposições a respeito da forma da distribuição.
- Custo computacional significativo para determinar a hipótese ótima de Bayes
 - É geralmente linear com o número de hipóteses

Teorema de Bayes

$P(c|X)$: probabilidade da classe c dado o vetor X

$P(X|c)$: probabilidade do vetor X dada a classe c .

$P(c)$: probabilidade *a priori* da classe c

$$P(c | X) = \frac{P(X | c) P(c)}{P(X)}$$

$P(X)$: probabilidade *a priori* do vetor de treinamento X

Teorema de *Bayes*

- $P(c|X)$ é chamada de probabilidade *a posteriori* de c porque ela reflete nossa confiança que c se mantenha após termos observado o vetor de treinamento X .
- $P(c|X)$ reflete a influência do vetor de treinamento X .
- Em contraste, a probabilidade *a priori* $P(c)$ é independente de X .

Teorema de *Bayes*

- Geralmente queremos encontrar a classe mais provável $c \in C$, sendo fornecidos os exemplos de treinamento X .
- Ou seja, a classe com o máximo *a posteriori* (*MAP*)

$$\begin{aligned}c_{MAP} &\equiv \arg \max_{c \in C} P(c | X) \\&= \arg \max_{c \in C} \frac{P(X | c)P(c)}{P(X)} \\&= \arg \max_{c \in C} P(X | c)P(c)\end{aligned}$$

Teorema de *Bayes*

- Desprezamos o termo $P(X)$ porque ele é uma constante independente de c .
- Se assumirmos que cada classe em C é igualmente provável *a priori*, i.e.

$$P(c_i) = P(c_j) \quad \forall \quad c_i \text{ e } c_j \text{ em } C$$

- Então, podemos simplificar e escolher a classe de máxima probabilidade condicional (*maximum likelihood* = *ML*).

Teorema de *Bayes*

- O termo $P(X|c)$ é chamado de probabilidade condicional (ou *likelihood*) de X
- Sendo fornecido c , qualquer classe que maximiza $P(X|c)$ é chamada de uma hipótese ML.

$$c_{ML} \equiv \arg \max_{c \in C} P(X | c)$$

Teorema de Bayes: Exemplo

- Considere um problema de diagnóstico médico onde existem duas classes possíveis:
 - O paciente tem H1N1
 - O paciente não tem H1N1
- As características disponíveis são um exame de laboratório com dois resultados possíveis:
 - \oplus : positivo
 - \ominus : negativo

Teorema de Bayes: Exemplo

- Temos o conhecimento prévio que na população inteira somente 0,008 tem esta doença.
- O exame retorna um resultado positivo correto somente em 98% dos casos nos quais a doença está presente.
- O exame retorna um resultado negativo correto somente em 97% dos casos nos quais a doença não esteja presente.
- Nos outros casos, o teste retorna o resultado oposto.

Teorema de Bayes: Exemplo

- $P(H1N1) = ?$

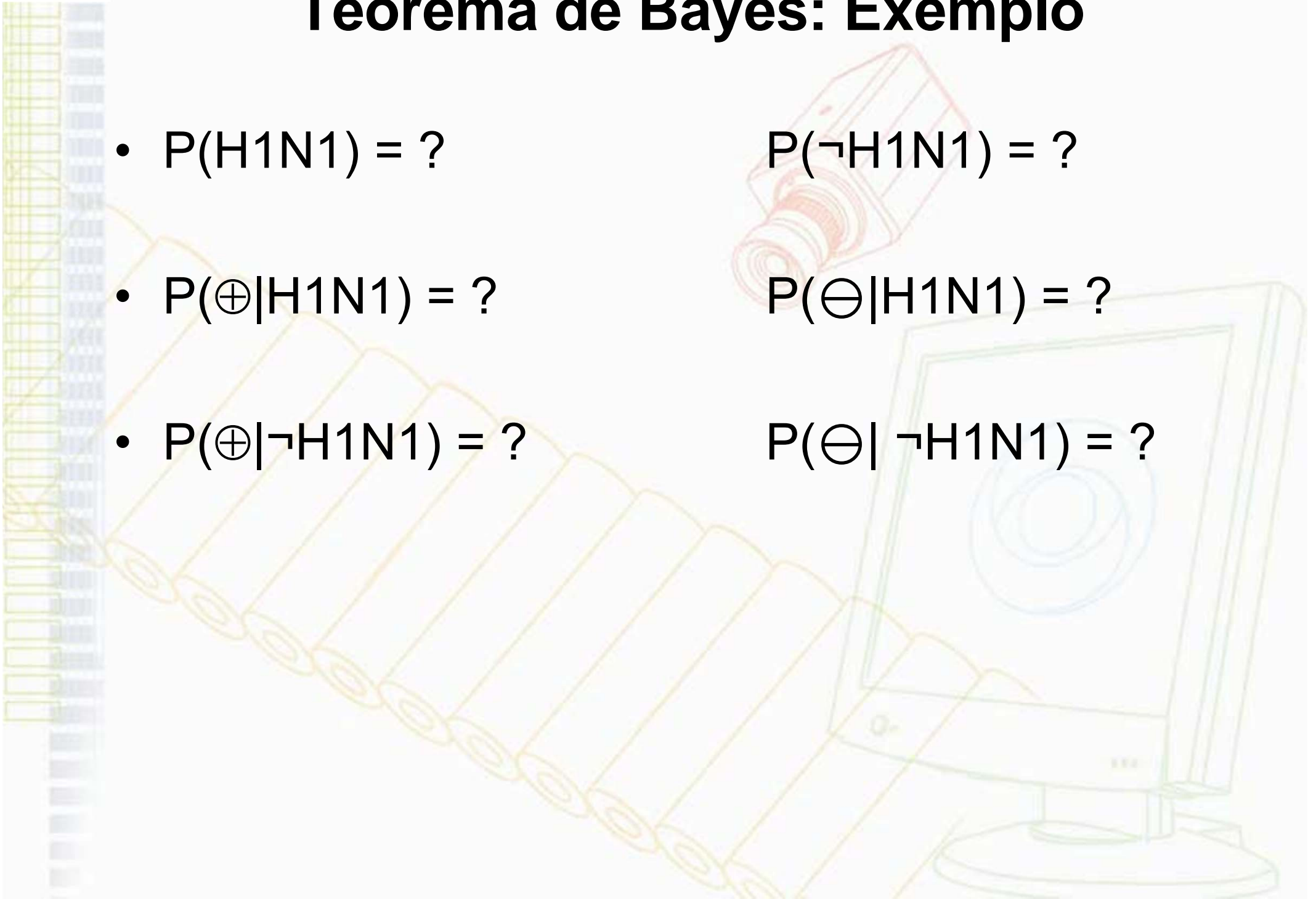
$$P(\neg H1N1) = ?$$

- $P(\oplus|H1N1) = ?$

$$P(\ominus|H1N1) = ?$$

- $P(\oplus|\neg H1N1) = ?$

$$P(\ominus|\neg H1N1) = ?$$



Teorema de Bayes: Exemplo

- Supondo que um paciente fez um exame de laboratório e o resultado deu positivo.
- O paciente tem H1N1 ou não ?

Aplicando o Teorema de Bayes

- Calculando a classe com maior probabilidade *a posteriori*:
 - $P(\oplus|H1N1) P(H1N1) = 0,98 \times 0,008 = 0,0078$
 - $P(\oplus|\neg H1N1) P(\neg H1N1) = 0,03 \times 0,992 = 0,0298$
- Assim: $c_{MAP} = \neg H1N1$

Formulação Básica de Probabilidades

- *Product rule*: probability $P(A \wedge B)$ of a conjunction of two events A and B

$$P(A \wedge B) = P(A|B)P(B) = P(B|A)P(A)$$

- *Sum rule*: probability of a disjunction of two events A and B

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

- *Bayes theorem*: the posterior probability $P(h|D)$ of h given D

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

- *Theorem of total probability*: if events A_1, \dots, A_n are mutually exclusive with $\sum_{i=1}^n P(A_i) = 1$, then

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

Classificador Ótimo de Bayes

- Consideramos até agora a questão:
“Qual a classe mais provável (c_{MAP}) dado os exemplos de treinamento X ?”
- Entretanto, a questão mais significativa é na verdade:
“Qual é a classificação mais provável de uma nova instância dado os dados de treinamento?”
- A classe MAP (c_{MAP}) é ou não a classificação mais provável?

Classificador Ótimo de Bayes

- Considere três classes possíveis c_1 , c_2 e c_3 e suponha as seguintes probabilidades *a posteriori* destas classes o conjunto de treinamento X :

$$P(c_1|X) = 0.4 \quad P(c_2|X) = 0.3 \quad P(c_3|X) = 0.3$$

- Qual é a classe *MAP*?

Classificador Ótimo de Bayes

- A classificação mais provável de uma nova instância x é obtida através da maior probabilidade *a posteriori*.
- Assim, a $P(c_j|x)$ que a correta classificação para a instância x seja c_j é:

$$\hat{P}(c_j | x) = \max_{c_j \in C} P(c_j | x)$$

$$\hat{c} = \arg \max_{c_j \in C} P(c_j | x)$$

- Qualquer sistema que classifique novas instâncias de acordo com a equação acima é chamada de um *classificador ótimo de Bayes*.

Exemplo

- Exemplo: Considere as 14 instâncias de treinamento de *PlayTennis* e uma nova instância de teste (x_t) que devemos classificar:

$x_t = \langle \text{Outlook}=\text{sunny}, \text{Temperature}=\text{cool}, \text{Humidity}=\text{high}, \text{Wind}=\text{strong} \rangle$

- Nossa tarefa é predizer o valor alvo (*yes* ou *no*) do conceito *PlayTennis* para esta nova instância, ou seja:

$$\hat{P}(c_j | x_t) = \max_{c_j \in [\text{yes}, \text{no}]} P(c_j | x_t)$$

$$\hat{c} = \arg \max_{c_j \in [\text{yes}, \text{no}]} P(c_j | x_t)$$

Exemplo

- Então, dado x_t , devemos estimar duas probabilidades *a posteriori*:

$$P(c_j = \text{yes} \mid x_t) \quad P(c_j = \text{no} \mid x_t)$$

- Aplicando o teorema de Bayes...

$$P(c_j \mid x_t) = \frac{P(x_t \mid c_j)P(c_j)}{P(x_t)}$$

- Ou seja, para estimar a probabilidade *a posteriori*, devemos conhecer:

$$P(x_t) = ?$$

$$P(x_t \mid c_j) = ?$$

$$P(c_j) = ?$$

Exemplo

- Atributo alvo: *PlayTennis* (yes, no)

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Exemplo

- Logo, temos que estimar:
 - duas probabilidades *a priori* das classes:
 $P(yes) = ? \quad P(no) = ?$
 - probabilidade *a priori* do vetor x_t :
 $P(x_t) = ?$
 - duas probabilidades condicionais:
 $P(x_t | yes) = ? \quad P(x_t | no) = ?$
- Como fazer isso dadas as 14 instâncias de treinamento da tabela?

Exemplo

$$P(yes) = 9 / 14 = 0,643$$

$$P(no) = 5 / 14 = 0,357$$

$$P(\langle outlook = sunny, temperature = hot, humidity = high, wind = weak \rangle | yes) = ?$$

$$P(\langle outlook = overcast, temperature = hot, humidity = high, wind = weak \rangle | yes) = ?$$

$$P(\langle outlook = rain, temperature = hot, humidity = high, wind = weak \rangle | yes) = ?$$

...

$$P(\langle outlook = rain, temperature = cool, humidity = normal, wind = strong \rangle | yes) = ?$$

....ou seja, temos que estimar todas as probabilidades condicionais, considerando todas as classes possíveis e todos os vetores de características possíveis:

$$2 \times [3 \times 3 \times 2 \times 2] = 72 \text{ probabilidades condicionais}$$

Exemplo

$$2 \times [3 \times 3 \times 2 \times 2] = 72$$

pois:

- temos 2 classes
- temos 4 atributos e seus possíveis valores:
 - Outlook (sunny/overcast/rain) [3 valores possíveis]
 - Temperature (hot/mild/cool) [3 valores possíveis]
 - Humidity (high/normal) [2 valores possíveis]
 - Wind (weak/strong) [2 valores possíveis]
- Logo, temos 72 probabilidades condicionais possíveis.
- e $P(x_t)$?

Classificador Ótimo de Bayes

- Limitações práticas
 - Como estimar com confiança todas estas probabilidades condicionais?
 - Conjunto de treinamento com muitas instâncias!
 - Conhecer a distribuição de probabilidade!
 - A probabilidade *a priori* calculada geralmente não reflete a população.

Classificador Naïve Bayes

- *Naïve Bayes* é um dos métodos de aprendizagem mais práticos.
- Quando usar ?
 - disponibilidade de um conjunto de treinamento grande ou moderado.
 - os atributos que descrevem as instâncias forem **condicionalmente independentes** dada a classe.
- Aplicações bem sucedidas:
 - diagnóstico médico
 - classificação de documentos de textuais

Classificador Naïve Bayes

- Se aplica a tarefas de aprendizagem onde:
 - cada instância x é descrita por uma conjunção de valores de atributos
 - a função alvo $f(x)$ pode assumir qualquer valor de um conjunto V .
 - um conjunto de exemplos de treinamento da função alvo é fornecido
 - uma nova instância é descrita pela *tupla* de valores de atributos $\langle a_1, a_2, \dots, a_n \rangle$.
- A tarefa é prever o valor alvo (ou classe) para esta nova instância.

Classificador Naïve Bayes

- A solução *Bayesiana* para classificar uma nova instância consiste em:
 - atribuir o valor alvo mais provável (c_{MAP}) dados os valores dos atributos $\langle a_1, a_2, \dots, a_n \rangle$ que descrevem a instância.

$$c_{MAP} = \arg \max_{c_j \in C} P(c_j | a_1, a_2, \dots, a_n)$$

- Mas podemos usar o teorema de Bayes para reescrever a expressão . . .

Classificador Naïve Bayes

$$c_{MAP} = \arg \max_{c_j \in C} P(c_j | a_1, a_2, \dots, a_n)$$

$$c_{MAP} = \arg \max_{c_j \in C} \frac{P(a_1, a_2, \dots, a_n | c_j) P(c_j)}{P(a_1, a_2, \dots, a_n)}$$

$$= \arg \max_{c_j \in C} P(a_1, a_2, \dots, a_n | c_j) P(c_j)$$

- Devemos agora estimar os dois termos da equação acima baseando-se nos dados de treinamento.
 - $P(c_j)$ é fácil de estimar . . .
 - Porém, $P(a_1, a_2, \dots, a_n | c_j)$. . .

Classificador Naïve Bayes

- O classificador Naïve Bayes é baseado na suposição simplificadora de que os valores dos atributos são condicionalmente independentes dado o valor alvo.
- Ou seja, a probabilidade de observar a conjunção de atributos a_1, a_2, \dots, a_n é somente o produto das probabilidades para os atributos individuais:

$$P(a_1, a_2, \dots, a_n | c_j) = \prod_i P(a_i | c_j)$$

Classificador Naïve Bayes

- Temos assim o classificador Naïve Bayes:

$$\hat{c}_{NB} = \arg \max_{c_j \in C} P(c_j) \prod_i P(a_i | c_j)$$

onde c_{NB} indica o valor alvo fornecido pelo algoritmo Naïve Bayes.

Classificador Naïve Bayes

- Em resumo, o algoritmo Naïve Bayes envolve
 - Aprendizagem: os termos $P(c_j)$ e $P(a_i|c_j)$ são estimados baseado nas suas frequências no conjunto de treinamento.
 - Estas probabilidades “aprendidas” são então utilizadas para classificar uma nova instância aplicando a equação vista anteriormente (c_{NB})

Classificador Naïve Bayes

Algoritmo Naïve Bayes

Treinamento_Naïve_Bayes(*conjunto de exemplos*)

Para cada valor alvo (classe) c_j

$P'(c_j) \leftarrow$ estimar $P(c_j)$

Para cada valor de atributo a_i de cada atributo a

$P'(a_i | c_j) \leftarrow$ estimar $P(a_i | c_j)$

Classifica_Naïve_Bayes(x_t)

$$\hat{c}_{NB} = \arg \max_{c_j \in C} P'(c_j) \prod_{a_i \in x} P'(a_i | c_j)$$

Classificador Naïve Bayes

- Exemplo: Considere novamente os 14 exemplos de treinamento de *PlayTennis* e uma nova instância que o Naïve Bayes deve classificar:

$x_t = \langle \text{outlook}=\text{sunny}, \text{temperature}=\text{cool}, \text{humidity}=\text{high}, \text{wind}=\text{strong} \rangle$

- A tarefa é predizer o valor alvo (*yes* ou *no*) do conceito *PlayTennis* para esta nova instância.

Classificador Naïve Bayes

- Atributo alvo: *PlayTennis* (yes, no)

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Classificador Naïve Bayes

- O valor alvo c_{NB} será dado por:

$$\begin{aligned} c_{NB} &= \arg \max_{c_j \in \{yes, no\}} P(c_j) \prod_i P(a_i | c_j) \\ &= \arg \max_{c_j \in \{yes, no\}} P(c_j) P(Outlook = sunny | c_j) P(Temperature = cool | c_j) \\ &\quad P(Humidity = high | c_j) P(Wind = strong | c_j) \end{aligned}$$

- Note que a_i foi instanciado utilizando os valores particulares do atributo da instância x_t .
- Para calcular c_{NB} são necessárias 10 probabilidades que podem ser estimadas a partir dos exemplos de treinamento.

Classificador Naïve Bayes

- Probabilidades *a priori*:

$$P(\text{PlayTennis} = \text{yes}) = 9/14 = 0.64$$

$$P(\text{PlayTennis} = \text{no}) = 5/14 = 0.36$$

- Probabilidades condicionais:

$$P(\text{Wind}=\text{strong} \mid \text{PlayTennis} = \text{yes}) = 3/9 = 0.33$$

$$P(\text{Wind}=\text{strong} \mid \text{PlayTennis} = \text{no}) = 3/5 = 0.60$$

...

Classificador Naïve Bayes

- Usando estas estimativas de probabilidade e estimativas similares para os valores restantes dos atributos, calculamos c_{NB} de acordo com a equação anterior (omitindo nome dos atributos) :

$$P(\text{yes}) P(\text{sunny} | \text{yes}) P(\text{cool} | \text{yes}) P(\text{high} | \text{yes}) P(\text{strong} | \text{yes}) = 0,0053$$

$$P(\text{no}) P(\text{sunny} | \text{no}) P(\text{cool} | \text{no}) P(\text{high} | \text{no}) P(\text{strong} | \text{no}) = 0,026$$

- Então o classificador atribui o valor alvo $\text{PlayTennis} = \text{no}$ para esta nova instância.

Classificador Naïve Bayes

Sutilezas:

1. A suposição de independência condicional é muitas vezes violada

$$P(a_1, a_2, \dots, c_j) = \prod_i P(a_i | c_j)$$

... mas, de qualquer maneira, ele funciona bem. Note que não é necessário estimar probabilidades *a posteriori* $P'(c_j|x)$ para ser correta. Necessita somente que

$$\arg \max_{c_j \in C} P'(c_j) \prod_i P'(a_i | c_j) = \arg \max_{c_j \in C} P(c_j) P(a_1, \dots, a_n | c_j)$$

- Probabilidades Naïve Bayes *a posteriori* próximas de 0 e 1 são geralmente não realísticas

Classificador Naïve Bayes

Sutilezas:

2. E se nenhuma das instâncias de treinamento com valor alvo c_j tiver um atributo de valor a_i ? Então,

$$P'(a_i | c_j) = 0$$

e ...

$$P'(c_j) \prod_i P'(a_i | c_j) = 0$$

A solução típica é uma estimativa Bayesiana para $P'(a_i | c_j)$

$$P'(a_i | c_i) \leftarrow \frac{n_c + mp}{n + m}$$

Classificador Naïve Bayes

$$P'(a_i | c_i) \leftarrow \frac{n_c + mp}{n + m}$$

onde:

- n é o número de exemplos de treinamento para os quais $c = c_j$,
- n_c é o número de exemplos para os quais $c = c_j$ e $a = a_i$
- p é a estimativa a priori para $P'(a_i | c_j)$
- m é o peso dado *as priori* (i.e. número de exemplos “virtuais”).

Exemplo: Classificando Texto

- Por que ?
 - Aprender quais notícias são interessantes
 - Aprender a classificar páginas WEB por assunto
 - Naïve Bayes é um dos algoritmos mais eficientes
 - Quais atributos devemos usar para representar documentos de texto?

Exemplo: Classificando Texto

- Contexto
 - Considere um espaço de instâncias X consistindo de todos os documentos de texto possíveis.
 - Dados exemplos de treinamento, de alguma função alvo $f(x)$ que pode assumir valores de um conjunto finito C .
 - A tarefa de aprendizagem é aprender, a partir dos exemplos de treinamento, a prever o valor alvo para os documentos de texto subsequentes.
 - Considere a função alvo como sendo documentos interessantes e não interessantes

Exemplo: Classificando Texto

- Projeto do Naïve Bayes:
 - Como representar um documento de texto arbitrário em termos de valores de atributos?
 - Decidir como estimar as probabilidades necessárias para o Naïve Bayes.

Exemplo: Classificando Texto

- Representação de texto arbitrário
 - Dado um documento de texto, este parágrafo, por exemplo, definimos um atributo para cada posição de palavra no documento e definimos o valor do atributo como sendo a palavra em português encontrada nesta posição.
 - O parágrafo anterior pode ser descrito por 34 valores de atributos correspondendo as 34 posições de palavras.
 - O valor do primeiro atributo é a palavra “Dado” e do segundo é a palavra “um” e assim por diante.

Exemplo: Classificando Texto

- Dada a representação de documento de texto, podemos aplicar o Naïve Bayes.
- Assumimos
 - um conjunto de 700 documentos classificados por uma pessoa como não interessantes
 - outros 300 classificados como interessantes

Exemplo: Classificando Texto

- Conceito alvo interessante: documento $\rightarrow \{+, -\}$
- 1. Representar cada documento por um vetor de palavras
 - Um atributo por posição da palavra no documento
- 2. Aprendendo usar exemplos de treinamento para estimar
 - $P(+)$
 - $P(-)$
 - $P(doc|+)$
 - $P(doc|-)$

Exemplo: Classificando Texto

- Suposição da independência condicional Naïve Bayes

$$P(doc | c_j) = \prod_{i=1}^{length(doc)} P(a_i = w_k | c_j)$$

onde $P(a_i = w_k | c_j)$ é a probabilidade que a palavra na posição i é w_k , dado c_j .

- Mais uma suposição

$$P(a_i = w_k | c_j) = P(a_m = w_k | c_j) \quad \forall i, m$$

Exemplo: Classificando Texto

Learn_Naïve_Bayes_Text (Examples, C)

1. Colectionar todas palavras, pontuação e outros tokens que ocorrem em *Examples*
 - *Vocabulary* \leftarrow todas as palavras distintas e outros tokens que ocorrem em *Examples*
2. Calcular as probabilidade necessárias $P(c_j)$ e $P(w_k|c_j) \dots$

Exemplo: Classificando Texto

- Para cada valor alvo c_j em V faça
 - $docs_j \leftarrow$ subconjunto de documento de *Examples* para o qual o valor alvo é c_j
 - $P(c_j) = \frac{|docs_j|}{|Examples|}$
 - $Text_j \leftarrow$ um documento único criado pela concatenação de todos os membros de $docs_j$
 - $n \leftarrow$ número total de posições distintas de palavras em $Text_j$
 - Para cada palavra w_k em *Vocabulary*
 - $n_k \leftarrow$ número de vezes que a palavra w_k ocorre em $Text_j$
 - $P(w_k | c_j) \leftarrow \frac{n_k + 1}{n + |Vocabulary|}$

Exemplo: Classificando Texto

Classify_Naïve_Bayes_Text (*Doc*)

- *positions* \leftarrow todas as posições das palavras em *Doc* que contém tokens encontrados em *Vocabulary*
- retornar c_{NB} onde

$$c_{NB} = \arg \max_{c_j \in C} P(c_j) \prod_{i \in positions} P(a_i | c_j)$$

Exemplo: Classificando Texto

LEARN_NAIVE_BAYES_TEXT(*Examples*, *V*)

Examples is a set of text documents along with their target values. *V* is the set of all possible target values. This function learns the probability terms $P(w_k|v_j)$, describing the probability that a randomly drawn word from a document in class v_j will be the English word w_k . It also learns the class prior probabilities $P(v_j)$.

1. collect all words, punctuation, and other tokens that occur in *Examples*
 - *Vocabulary* \leftarrow the set of all distinct words and other tokens occurring in any text document from *Examples*
2. calculate the required $P(v_j)$ and $P(w_k|v_j)$ probability terms
 - For each target value v_j in *V* do
 - *docs_j* \leftarrow the subset of documents from *Examples* for which the target value is v_j
 - $P(v_j) \leftarrow \frac{|docs_j|}{|Examples|}$
 - *Text_j* \leftarrow a single document created by concatenating all members of *docs_j*
 - *n* \leftarrow total number of distinct word positions in *Text_j*
 - for each word w_k in *Vocabulary*
 - $n_k \leftarrow$ number of times word w_k occurs in *Text_j*
 - $P(w_k|v_j) \leftarrow \frac{n_k+1}{n+|Vocabulary|}$

CLASSIFY_NAIVE_BAYES_TEXT(*Doc*)

Return the estimated target value for the document *Doc*. a_i denotes the word found in the *i*th position within *Doc*.

- *positions* \leftarrow all word positions in *Doc* that contain tokens found in *Vocabulary*
- Return v_{NB} , where

$$v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_{i \in \text{positions}} P(a_i|v_j)$$

TABLE 6.2

Naive Bayes algorithms for learning and classifying text. In addition to the usual naive Bayes assumptions, these algorithms assume the probability of a word occurring is independent of its position within the text.

Exemplo: Classificando Texto

- Dados 1.000 documentos de treinamento de cada grupo, aprenda a classificar novos documentos de acordo com o newsgroup de origem.

comp.graphics	misc.forsale	soc.religion.christian	sci.space
comp.os.ms-windows.misc	rec.autos	talk.politics.guns	sci.crypt
comp.sys.ibm.pc.hardware	rec.motorcycles	talk.politics.mideast	sci.electronics
comp.sys.mac.hardware	rec.sport.baseball	talk.politics.misc	sci.med
comp.windows.x	rec.sport.hockey	talk.religion.misc	
		alt.atheism	

TABLE 6.3

Twenty usenet newsgroups used in the text classification experiment. After training on 667 articles from each newsgroup, a naive Bayes classifier achieved an accuracy of 89% predicting to which newsgroup subsequent articles belonged. Random guessing would produce an accuracy of only 5%.

- Naïve Bayes: precisão de classificação: 89%

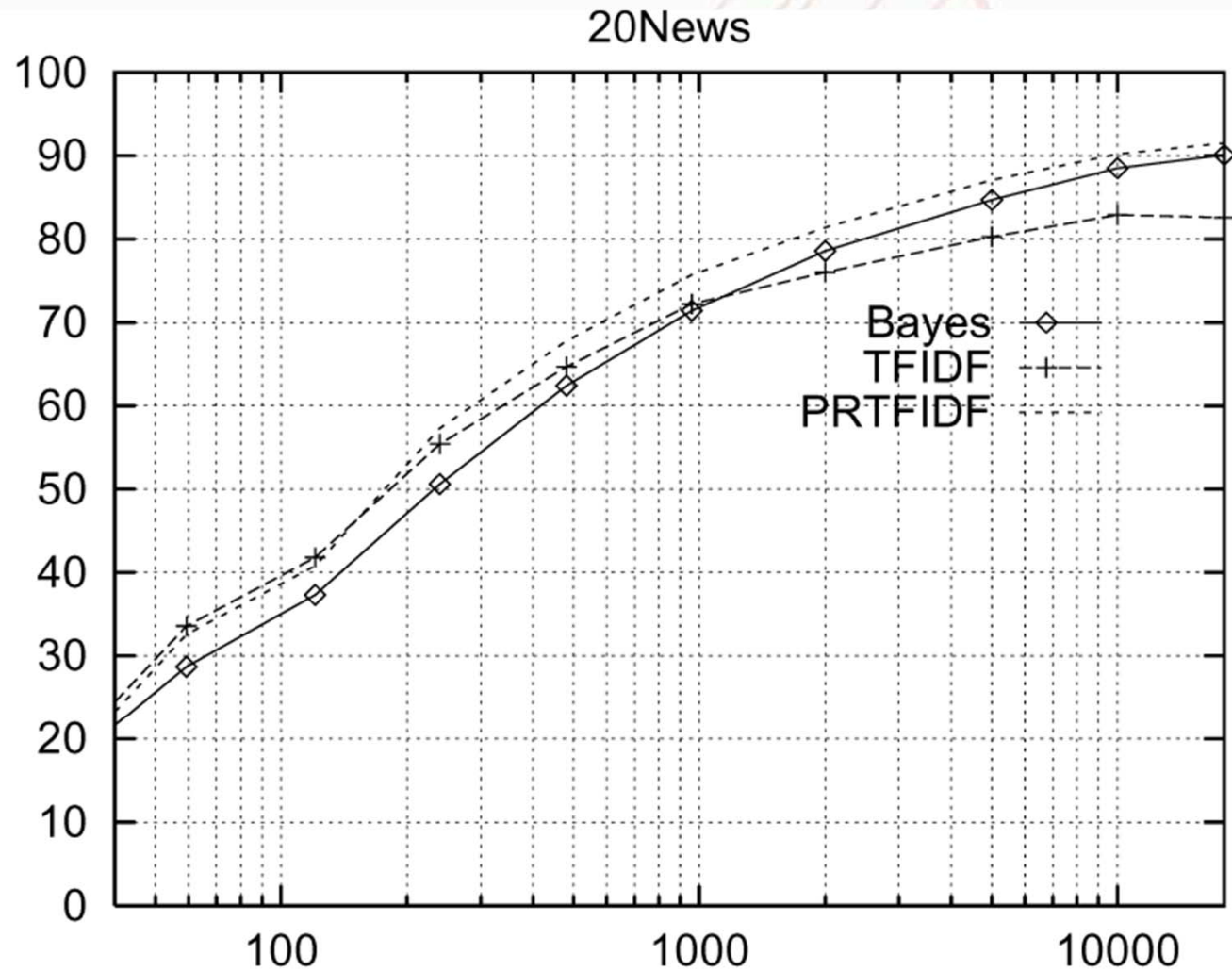
Exemplo: Classificando Texto

- Artigo de rec.sport.hockey

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e
From: xxx@yyy.zzz.edu (John Doe)
Subject: Re: This year's biggest and worst (opinio
Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrudey is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some thugs in Toronto decided

Curva de Aprendizagem



Accuracy vs. Training set size (1/3 withheld for test)

Resumo

- Métodos Bayesianos:
 - acomodam conhecimento prévio e os dados observáveis;
 - atribuem probabilidade a posteriori para cada classe candidata, baseando-se na probabilidade a priori e nos dados.
 - podem determinar a hipótese mais provável (MAP), tendo os dados.
- Bayes Ótimo:
 - combina previsões de todas classes, ponderadas pela probabilidade a posteriori, para calcular a classificação mais provável de uma nova instância.

Resumo

- Naïve Bayes:
 - é chamado de naïve (simples, não sofisticado), porque assume que os valores dos atributos são condicionalmente independentes.
 - se a condição é encontrada, ele fornece a classificação MAP, caso contrário, pode fornecer também bons resultados.