

LABIC

SISTEMAS INTELIGENTES

Profa. Roseli Ap. Francelin Romero

RAFR Sistemas Inteligentes 1

LABIC

Fórmulas Básicas para Probabilidades

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

$$P(A \wedge B) = P(A|B) P(B) = P(B|A) P(A)$$
- Regra Soma: probabilidade $P(A \vee B)$ de uma união de dois eventos **A** e **B**:

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$
- Teorema da probabilidade total: se eventos A_1, \dots, A_n são mutuamente exclusivos com $\sum_{i=1}^n P(A_i) = 1$, então:

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

RAFR Sistemas Inteligentes 2

LABIC

Aprendizado Bayesiano

CLASSIFICADORES BAYESIANO

Aprendizado Supervisionado de Classificadores Bayesiano

Aprendizado Não Supervisionado de Classificadores Bayesiano

RAFR Sistemas Inteligentes 3

LABIC

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?

RAFR Sistemas Inteligentes 4

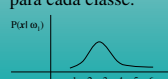
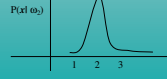
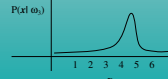
LABIC

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 *tirão* 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.

RAFR Sistemas Inteligentes 5

LABIC

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 ω_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).

RAFR Sistemas Inteligentes 6

LABIC

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 = $\arg \max P(\omega_i | x)$.

Bayesiano Ótimo = $\arg \max P(x | \omega_i) \cdot P(\omega_i)$

Este é o Classificador Ótimo de Bayes.

RAFR Sistemas Inteligentes 7

LABIC

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 = $\arg \max P(x | \omega_i) \cdot P(\omega_i)$

Hipótese Comum:

Cada $P(x | \omega_i)$ segue distribuição Gaussiana.

Três Casos:

- $P(x | \omega_i)$ - Média μ_i , variância σ^2
- $P(x | \omega_i)$ - Média μ_i , covariância Σ , arbitrária
- $P(x | \omega_i)$ - Média μ_i , covariância Σ_i , diferente para classes diferentes

RAFR Sistemas Inteligentes 8

LABIC

Caso 1

Caso 1: Todas componentes são independentes

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

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

Bayesiano Ótimo = $\arg \max P(x | \omega_i) \cdot P(\omega_i) = \arg \max \{ k \exp \left(\frac{-1}{2\sigma^2} \sum_j (x_j - \mu_{ij})^2 \right) \cdot P(\omega_i) \} =$

RAFR Sistemas Inteligentes 9

LABIC

Caso 1

$$= \arg \max \frac{-1}{2\sigma^2} \sum_j (x_j - \mu_{ij})^2 + \log P(\omega_i) =$$

$$= \arg \min \frac{\sum_j (x_j - \mu_{ij})^2 - 2\sigma^2 \log P(\omega_i)}{2\sigma^2} =$$

$$= \arg \min \sum_j (x_j - \mu_{ij})^2 - 2\sigma^2 \log P(\omega_i)$$

RAFR Sistemas Inteligentes 10

LABIC

Caso 1

■ Caso duas classes

$$= \arg \min (\sum_j (x_j - \mu_1)^2 - 2\sigma^2 \log P(\omega_1)) =$$

$$= \arg \min (\sum_j x_j x_j - 2 \sum_j x_j \mu_1 + \sum_j \mu_1 \mu_1 - 2\sigma^2 \log P(\omega_1)) =$$

$$= \arg \min (\sum_j - 2 \sum_j x_j \mu_1 + c_1)$$

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

RAFR Sistemas Inteligentes 11

LABIC

Caso 1

\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.

RAFR Sistemas Inteligentes 12

LABIC

Caso 2 - Hipótese mais fraca

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

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

Se todos os $P(\omega_i)_s$ são iguais \Rightarrow método do vizinho mais próximo (KNN)
 Ainda usa regiões de decisão linear.

RAFR Sistemas Inteligentes 13

LABIC

Caso 3 - Hipótese ainda mais fraca

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

Ainda é fácil calcular a decisão ótima

$$\text{arg,max}(P(\omega_i | x))$$

mas as regiões de decisão não são mais lineares.

RAFR Sistemas Inteligentes 14

LABIC

Classificacao de Padroes

- Suponha agora que voce não 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.

$x_1^{(1)} x_2^{(1)} \dots x_N^{(1)}$	Classe w_1
$x_1^{(2)} x_2^{(2)} \dots x_N^{(2)}$	Classe w_2
..	
$x_1^{(M)} x_2^{(M)} \dots x_N^{(M)}$	Classe w_N

RAFR Sistemas Inteligentes 15

LABIC

Classificacao de Padroes

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

$$\text{arg,max } P(x|\omega_i).P(\omega_i)$$

RAFR Sistemas Inteligentes 16

LABIC

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.

RAFR Sistemas Inteligentes 17

LABIC

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.

RAFR Sistemas Inteligentes 18

LABIC

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.

RAFR Sistemas Inteligentes 19

LABIC

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.

RAFR Sistemas Inteligentes 20

LABIC

TEOREMA DE BAYES

Em problemas de ML estamos interessados em $P(h|D)$: probabilidade a posteriori, probabilidade que 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 uma hipótese mais provável, se existirem várias).

RAFR Sistemas Inteligentes 21

LABIC

TEOREMA DE BAYES

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

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

É independente de $P(D)$

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:

RAFR Sistemas Inteligentes 22

LABIC

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 - o espaço das funções alvo candidatas.

RAFR Sistemas Inteligentes 23

LABIC

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.

$P(\text{câncer}) = 0.008$ $P(\neg \text{câncer}) = 0.992$
 $P(+|\text{câncer}) = 0.98$ $P(-|\text{câncer}) = 0.02$
 $P(+|\neg \text{câncer}) = 0.03$ $P(-|\neg \text{câncer}) = 0.97$
 $P(+|\text{câncer}) \cdot P(\text{câncer}) = (0.98) \cdot (0.008) = 0.0078$
 $P(+|\neg \text{câncer}) \cdot P(\neg \text{câncer}) = (0.03) \cdot (0.992) = 0.0298$

$h_{MAP} = \neg \text{câncer}$

RAFR Sistemas Inteligentes 24

Classificação mais Provável de Novas Instâncias

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

Dado nova instância x , qual é a sua classificação mais provável?

$h_{MAP}(x)$ 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 ?

$p_+(x)=0.4, p_-(x)=0.6$, portanto é mais provável que x seja -

Neste caso, é diferente da classificação gerada pela h_{MAP}

LABIC RAFR Sistemas Inteligentes 25

Classificador Bayesiano Ótimo

$$\arg_{v_j \in V} \max_{h \in H} \sum P(v_j | h_i) \cdot P(h_i | D)$$

EXEMPLO:

$P(h_1 | D) = 0.4, P(- | h_1) = 0, P(+ | h_1) = 1,$
 $P(h_2 | D) = 0.3, P(- | h_2) = 1, P(+ | h_2) = 0,$
 $P(h_3 | D) = 0.3, P(- | h_3) = 1, P(+ | h_3) = 0,$

Portanto, $\sum_{h \in H} P(+ | h_i) \cdot P(h_i | D) = 0.4$
 $\sum_{h \in H} P(- | h_i) \cdot P(h_i | D) = 0.6$

Portanto, $\arg_{v_j \in V} \max_{h \in H} \sum P(v_j | h_i) \cdot P(h_i | D) = -$

LABIC RAFR Sistemas Inteligentes 26

Aprendizado de uma Função Real

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

$$d_i = f(x_i) + e_i$$

onde e_i é uma variável aleatória, independente para

LABIC RAFR Sistemas Inteligentes 27

Aprendizado de uma Função Real

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

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

Demonstração:

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

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

Maximizando o logaritmo natural:

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

LABIC RAFR Sistemas Inteligentes 28

Aprendizado de uma Função Real

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

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

LABIC RAFR Sistemas Inteligentes 29

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

LABIC RAFR Sistemas Inteligentes 30

LABIC

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

RAFR Sistemas Inteligentes 31

LABIC

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:
 $\langle \text{Outlook} = \text{sunny}, \text{Temp} = \text{cool}, \text{Hum} = \text{high}, \text{wind} = \text{strong} \rangle$
 Queremos:

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

RAFR Sistemas Inteligentes 32

LABIC

Algoritmo Naive Bayes

- Naive_Bayes_Learn(examples)
 - Para cada valor alvo v_j
 - Calcule $P'(v_j) \leftarrow$ estimativa $P(v_j)$
 - Para cada valor de atributo a_i de cada atrib. A
 - $P'(a_i | v_j) \leftarrow$ estimativa $P(a_i | v_j)$

Classificar nova instancia x

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

RAFR Sistemas Inteligentes 33

LABIC

Exemplos de Treinamento

DAY	OUTLOOK	TEMPERATURA	UMIDADE	VENTO	PLAYTENN
D1	SOL	QUENTE	ALTA	FRACO	NÃO
D2	SOL	QUENTE	ALTA	FORTE	NÃO
D3	NUBLADO	QUENTE	ALTA	FRACO	SIM
D4	CHUVA	AMENO	ALTA	FRACO	SIM
D5	CHUVA	FRIO	NORMAL	FRACO	SIM
D6	CHUVA	FRIO	NORMAL	FORTE	NÃO
D7	NUBLADO	FRIO	NORMAL	FORTE	SIM
D8	SOL	AMENO	ALTA	FRACO	NÃO
D9	SOL	FRIO	NORMAL	FRACO	SIM
D10	CHUVA	AMENO	NORMAL	FRACO	SIM
D11	SOL	AMENO	NORMAL	FORTE	SIM
D12	NUBLADO	AMENO	ALTA	FORTE	SIM
D13	NUBLADO	QUENTE	NORMAL	FRACO	SIM
D14	CHUVA	AMENO	ALTA	FORTE	NÃO

RAFR Sistemas Inteligentes 34

LABIC

EXEMPLO

Considere o exemplo "Play Tennis" e a instância:
 $\langle \text{Outlook} = \text{sunny}, \text{Temp} = \text{cool}, \text{Hum} = \text{high}, \text{wind} = \text{strong} \rangle$
 Queremos:

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

RAFR Sistemas Inteligentes 35

LABIC

Classificador Bayesiano Naive

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

$$9/14 * 2/9 * 3/9 * 3/9 * 3/9 =$$

$$= 0.0053$$

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

$$= 0.0206$$

$\rightarrow V_{NB} = n$

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

RAFR Sistemas Inteligentes 36



Curiosidades

- A hipótese de independência condicional é frequentemente violada
Mas, mesmo assim o Classificador Naive Bayes funciona muito bem.
- Não se precisa de estimativas corretas da prob. A posteriori $P(v_j | x)$ mas apenas de:
 $P(v_j)$ e $P(a_i | v_j)$



Curiosidades

- O que fazer se nenhum dos exemplos no conj. De dados com valor alvo v_j tem valor de atributo a_i ?

$$P'(a_i | v_j) = 0 \quad \text{e}$$

$$P'(v_j) \pi P'(a_i | v_j) = 0$$

Solução é estimar $P'(a_i | v_j)$

$$P'(a_i | v_j) = (n_c + m) / (n + m \cdot p)$$

Onde: n - no. exs. de treinamento para o qual $v = v_j$

n_c - no. de exs. Virtuais para os quais $v = v_j$ e $a = a_i$

p é uma estimativa da priori para $P'(a_i | v_j)$ e m é um peso dado a priori