

Aprendizado de Máquina: Aprendizado Bayesiano

Prof. Arnaldo Candido Junior
UTFPR – Medianeira

Revisão de probabilidades

- Considere a tabela ao lado
- Uma pessoa nasce com olho azul quando possui dois genes recessivos

Mãe	Pai	Filho
a	a	aa
a	C	aC
C	a	Ca
C	C	CC

- Para simplificar:
 - Só serão analisados olhos azuis e castanhos
 - Será considerado que os genes estão igualmente distribuídos na população

Revisão de probabilidades (2)

- Revisão:
 - Probabilidade a priori
 - Probabilidade conjunta
 - Probabilidade a posteriori
- Teorema de Bayes
- Hipótese ingênua

Revisão de probabilidades ⁽³⁾

- Probabilidade **a priori**: calculada a partir daquilo que se sabe no início do problema
- Calcule as probabilidades dos eventos
 - A: filho tem olhos azuis
 - B: filho tem olhos castanhos
 - C: mãe forneceu gene recessivo
 - D: mãe forneceu gene dominante

Revisão de probabilidades ₍₄₎

- Para isso, é só contar as linhas da tabela que casam com os eventos (abordagem frequentista)
- $p(A) = 25\%$
- $p(B) = 75\%$
- $p(C) = 50\%$
- $p(D) = 50\%$

Revisão de probabilidades (5)

- Probabilidade **conjunta**: a chance de dois (ou mais) eventos ocorrem juntos
- Calcule
 - $p(A \wedge C)$
 - $p(B \wedge C)$
 - $p(B \wedge D)$

Revisão de probabilidades ₍₆₎

- Basta contar as linhas da tabela que casam com ambos eventos
 - $p(A \wedge C) = 25\%$
 - $p(B \wedge C) = 25\%$
 - $p(B \wedge D) = 50\%$

Revisão de probabilidades ₍₇₎

- Probabilidade **a posteriori**: calculada quando se recebe uma informação nova do problema
- Calcule
 - $p(A|C)$, lê-se $p(A)$ dado C
 - $p(B|D)$

Revisão de probabilidades ₍₈₎

- Por contagem: filtre pelo evento dado e conte o evento de interesse
 - $p(A|C) = 50\%$
 - $p(B|D) = 100\%$
- Por fórmula:

$$p(A|B) = \frac{p(A \wedge B)}{p(B)}$$

Teorema de Bayes

- Assuma $p(A) \neq 0$ e $p(B) \neq 0$

$$p(A|B) = \frac{p(A \wedge B)}{p(B)} \rightarrow p(A \wedge B) = p(A|B) p(B)$$

$$p(B|A) = \frac{p(A \wedge B)}{p(A)} \rightarrow p(A \wedge B) = p(B|A) p(A)$$

- Combinando as duas equações:

$$p(A|B) p(B) = p(B|A) p(A) \rightarrow p(A|B) = \frac{p(B|A) p(A)}{p(B)}$$

Teorema de Bayes ₍₃₎

- Exercício: para a tabela da cor dos olhos:
 - Calcule pelo método de contagem:
 $p(B)$
 $p(D)$
 $p(B|D)$
 $p(D|B)$
 - Calcule pelo teorema: $p(D|B)$

Teorema de Bayes ⁽⁴⁾

- Abordagem frequentista: consiste em contar linhas na tabela
- Abordagem bayesiana: consiste em usar aproximações para montar a tabela
 - Não dá para ter “clones” de um paciente para saber quantos dos clones tem determinada doença a partir dos sintomas
 - Usa-se outros pacientes com sintomas parecidos como uma aproximação do paciente para saber a doença

Hipótese Ingênua

- Hipótese Ingênua: $p(A \wedge B) = p(A)P(B)$
- Válida apenas quando A e B são dois eventos independentes
- Nem sempre é verdade: peso é completamente independente da altura?
- Exercício: use a hipótese ingênua para calcular $p(A \wedge C)$ na tabela de genética

Classificador Bayesiano Ingênuo

- Seja $f: X \rightarrow Y$ um problema de classificação
- Seja \hat{x} uma instância deste problema:
 $\hat{x} = \{x_1, \dots, x_n\}, \hat{x} \in X$
- Sejam v_i os valores dos atributos
- Desejamos calcular:

$$f = \underset{y_j \in Y}{\operatorname{argmax}} p(y = y_j | x_1 = v_1 \wedge x_2 = v_2 \wedge \dots \wedge x_n = v_n)$$

Classificador Bayesiano Ingênuo ⁽²⁾

- Ou seja, desejamos encontrar a classe mais provável da instância
- Primeiro, vamos simplificar a notação:

$$f = \underset{y_j}{\operatorname{argmax}} p(y_j | v_1 \wedge v_2 \wedge \dots \wedge v_n)$$

- Pelo **teorema de bayes**

$$f = \underset{y_i}{\operatorname{argmax}} \frac{p(v_1 \wedge v_2 \wedge \dots \wedge v_n | y_j) p(y_j)}{p(v_1 \wedge v_2 \wedge \dots \wedge v_n)}$$

Classificador Bayesiano Ingênuo ⁽³⁾

- Para analisar uma dada instância, o denominador é constante, podemos simplificar o cálculo para:

$$f = \underset{y_j}{\operatorname{argmax}} p(v_1 \wedge v_2 \wedge \dots \wedge v_n | y_j) p(y_j)$$

- Pela **hipótese ingênua**:

$$f = \underset{y_j}{\operatorname{argmax}} p(v_1 | y_j) p(v_2 | y_j) \dots p(v_n | y_j) p(y_j)$$

$$f = \underset{y_j}{\operatorname{argmax}} p(y_j) \prod_{v_i} p(v_i | y_j)$$

Classificador Bayesiano Ingênuo ⁽⁴⁾

- Estimativa das Probabilidades $p(v_i | y_j)$ e $p(y_j)$ pode
 - $p'(v_i | y_j) \leftarrow$ estimativa de $p(v_i | y_j)$
 - $p'(y_j) \leftarrow$ estimativa de $p(y_j)$
- Classificador de novas instancias(x)
$$h^* = \underset{y_j}{\operatorname{argmax}} p'(y_j) \prod_{v_i} p'(v_i | y_j)$$
- Idealmente, h^* aproxima f

Classificador Bayesiano Ingênuo (5)

•	<u>Dia</u>	<u>Tempo</u>	<u>Temp.</u>	<u>Humid.</u>	<u>Vento</u>	<u>Jogar</u>
•	D1	Sol	Quente	Alta	Fraco	Não
•	D2	Sol	Quente	Alta	Forte	Não
•	D3	Coberto	Quente	Alta	Fraco	Sim
•	D4	Chuva	Normal	Alta	Fraco	Sim
•	D5	Chuva	Frio	Normal	Fraco	Não
•	D6	Chuva	Frio	Normal	Forte	Não
•	D7	Coberto	Frio	Normal	Forte	Sim
•	D8	Sol	Normal	Alta	Fraco	Não
•	D9	Sol	Frio	Normal	Fraco	Sim
•	D10	Chuva	Normal	Normal	Fraco	Sim
•	D11	Sol	Frio	Alta	Forte	?

Classificador Bayesiano Ingênuo ⁽⁶⁾

$$p(\text{Sim}) = 5/10 = 0.5$$

$$p(\text{Não}) = 5/10 = 0.5$$

$$p(\text{Sol/Sim}) = 1/5 = 0.2$$

$$p(\text{Sol/Não}) = 3/5 = 0.6$$

$$p(\text{Frio/Sim}) = 2/5 = 0.4$$

$$p(\text{Frio/Não}) = 2/5 = 0.4$$

$$p(\text{Alta/Sim}) = 2/5 = 0.4$$

$$p(\text{Alta/Não}) = 3/5 = 0.6$$

$$p(\text{Forte/Sim}) = 1/5 = 0.2$$

$$p(\text{Forte/Não}) = 2/5 = 0.4$$

$$\begin{aligned} &P(\text{Sim})P(\text{Sol/Sim})P(\text{Frio/Sim}) \\ &P(\text{Alta/Sim})P(\text{Forte/Sim}) = \\ &= 0.0032 \end{aligned}$$

$$\begin{aligned} &P(\text{Não})P(\text{Sol/Não})P(\text{Frio/Não}) \\ &P(\text{Alta/Não})P(\text{Forte/Não}) = \\ &= 0.0288 \end{aligned}$$

⇒ Jogar_Tenis (D11) = Não

Exercício

- Faça um classificador bayesiano ingênuo para:

Nome	Febre	Enjôo	Manchas	Dores	Diagnóstico
João	sim	sim	pequenas	sim	doente
Pedro	não	não	grandes	não	saudável
Maria	sim	sim	pequenas	não	saudável
José	sim	não	grandes	sim	doente
Ana	sim	não	pequenas	sim	saudável
Leila	não	não	grandes	sim	doente

- (Luis, não, não, pequenas, sim)?
- (Laura, sim, sim, grandes, sim)?

Probabilidades nulas

- Novas instâncias podem conter valores de atributos não vistos em treinamento
 - Ex.: ceu = “nublado” na classe negativa não ocorreu
- Esses valores zeram os scores do produto para a nova instância
 - É possível que todas as possíveis classes fiquem com score 0

Probabilidades nulas ₍₂₎

- Solução intuitiva: adicione 1 em todos os nominadores

$$p(\text{Sol/Sim}) = 1/5 \rightarrow 2/5$$

$$p(\text{Sol/Não}) = 3/5 \rightarrow 4/5$$

$$p(\text{Frio/Sim}) = 2/5 \rightarrow 3/5$$

$$p(\text{Frio/Não}) = 2/5 \rightarrow 3/5$$

...

- Pode alterar o resultado da classe mais provável
- Solução utilizada: suavização de Laplace, de LidStone, entre outras

Suavização de Laplace

- Suavização de Laplace usa um pseudocontador
- No nominador, somar 1
- No denominador, somar o número de valores que um atributo pode assumir
- É uma forma de **data augmentation**: cria instâncias virtuais

Suavização de Laplace ₍₂₎

- Suponha que sol tem céu tenha 3 valores e que temperatura tenha 3

$$p(\text{Sol/Sim}) = 1/5 \rightarrow (1+1)/(5+3)$$

$$p(\text{Sol/Não}) = 3/5 \rightarrow (1+1)/(5+3)$$

$$p(\text{Frio/Sim}) = 2/5 \rightarrow (2+1)/(5+3)$$

$$p(\text{Frio/Não}) = 2/5 \rightarrow (3+1)/(5+3)$$

Outra Suavização: Text Mining

- Suponha que desejamos classificar textos em humorístico, esportivo ou biografia
- Os atributos representam as palavras que podem ocorrer no texto (bag of words)
- Os valores representam quantas vezes cada palavra ocorreu
- Problema: nem todas as palavras aparecem nos documentos de cada classe
- Isso leva a probabilidades nulas

Outra Suavização: Text Mining ₍₂₎

- Existem duas formas de calcular probabilidades da palavra x_i dada a classe y_j
- Opção 1: número de documentos da classe y_j que possuem a palavra x_i pelo número de total de documentos na classe y_j
- Opção 2: total de ocorrências da palavra x_i em documentos da classe y_j pelo total geral de palavras de todos os documentos da classe y_j

Outra Suavização: Text Mining ⁽³⁾

- O segundo caso ajuda a criar uma intuição sobre o processo de suavização.
 - Intuitivamente, colocamos as documentos de cada classe em uma pasta física
 - Depois, anexamos o dicionário Aurélio a cada pasta
 - Aí procedemos com a contagem
 - É garantido que cada palavra ocorrerá pelo menos uma vez em cada classe

Outra Suavização: Text Mining ₍₄₎

- $p'(y_j) = \frac{\text{número de documentos da classe } y_j}{\text{número total de documentos}}$
- $p'(v_i|y_j) = \frac{n_{ij} + 1}{n_j + |\text{Vocabulário}|}$
- n_j é o número **total** de palavras de **todos** os documentos da classe y_j
- n_{ij} é o número **total** de ocorrências da palavra x_i em **todos** os documentos da classe y_j

Log da probabilidade

- É comum usar logaritmo natural (**log-likelihood**) em implementações do classificador bayesiano ingênuo
 - Enquanto probabilidades são multiplicadas, logs são somados (mais eficiente)
 - Evita erros de arredondamento
 - Produtos de probabilidade convergem rapidamente a zero
 - A precisão dos tipos de ponto flutuante começa a deixar a deixar

Log da probabilidade ₍₂₎

- Revisão rápida de logaritmos (base 2):

- $8 * 32 = 256$

- $2^3 * 2^5 = 2^8$

- $\log(2^3 * 2^5) = \log 2^8$

- $\log 2^3 + \log 2^5 = \log 2^8$

- $3 + 5 = 8$

- $8 * 0.5 = 4$

- $2^3 * 2^x = 2^2$

- $x = -1$

Log da probabilidade ⁽³⁾

- Seja **p** uma probabilidade e **q** seu logaritmo natural

- Para recuperar p, é só fazer:

$$p = e^q$$

- Assim, a fórmula será revisada para:

$$h* = \underset{y_j}{\operatorname{argmax}} \left(\log p'(y_j) + \sum_{v_i} \log p'(v_i|y_j) \right)$$

Atributos contínuos

- Opção 1: discretizar
- Opção 2: considerar que os dados seguem uma distribuição normal
 - Usar a função de densidade de probabilidade normal para estimar a probabilidade

Redes Bayesianas

- Alternativa ao Bayesiano ingênuo
- Busca encontrar e incluir probabilidades dependentes no cálculo ao invés de recorrer à hipótese ingênua
- É um grafo acíclico e dirigido: cada nó da rede representa uma variável aleatória
 - Um conjunto de ligações ou arcos dirigidos conectam pares de nós

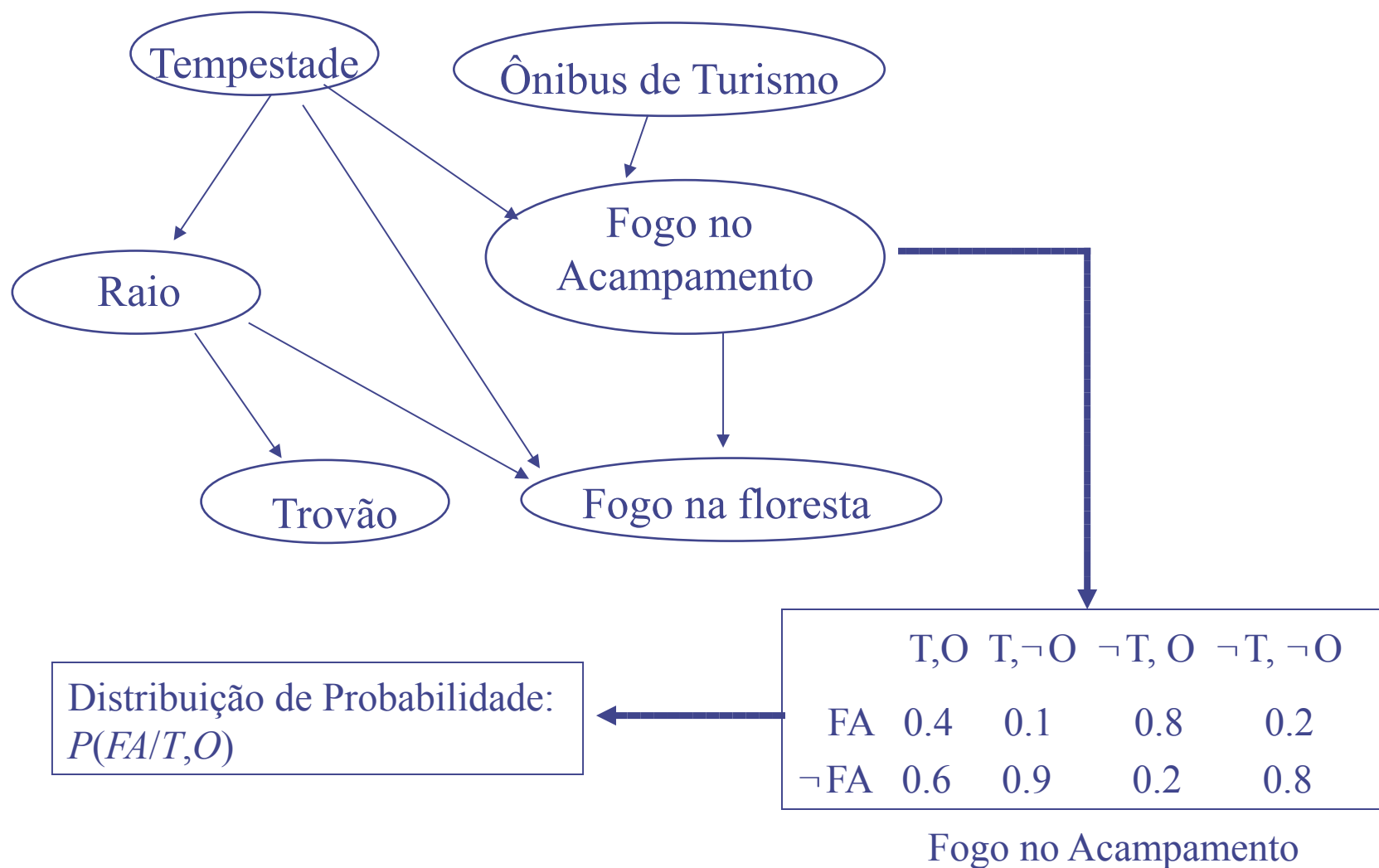
Redes Bayesianas ₍₂₎

- Cada nó recebe arcos dos nós que tem influência direta sobre ele (nós pais).
- Cada nó possui uma tabela de probabilidade condicional associada que quantifica os efeitos que os pais têm sobre ele
- Variações do algoritmo:
 - Projetista informa dependências entre atributos
 - Algoritmo encontra sozinho dependências entre atributos

Redes Bayesianas ⁽³⁾

- Como rede Bayesiana busca encontrar probabilidades conjuntas, é mais adequada para datasets maiores
 - Datasets grandes tem mais chances de conter probabilidades conjuntas
- Exemplo a seguir:
 - Rede para encontrar fogo na floresta
 - Para simplificar, pensar no caso fogo no acampamento

Redes Bayesianas (4)



Redes Bayesianas ₍₄₎

- No exemplo:
 - Estamos interessados em saber se há fogo no acampamento de posse das demais informações
 - Ônibus de turismo pode causar fogo no acampamento
 - Tempestade pode apagá-lo
 - Ou ainda, causar relâmpagos que causam fogo

Créditos

- Baseado em material da UFPE, disponível em www.cin.ufpe.br/~if684/aulas/Aprend_Bayestbl.ppt