

Naive Bayes

25 de Abril de 2015

Formulas de probabilidade básicas

$$P(A \wedge B) = P(A|B)P(B) = P(B|A)P(A) \quad (1)$$

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B) \quad (2)$$

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)} \quad (3)$$

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i) \text{ se } \sum_{i=1}^n P(A_i) = 1 \quad (4)$$

Classificador Naive Bayes

Seja a_1, \dots, a_n um conjunto de atributos e $v_i \in V$ um conjunto de valores de classificação possíveis. Então, de acordo com o teorema de *Bayes*, a probabilidade desse conjunto de atributos pertencer à classe v_i é

$$P(v_i | a_1, \dots, a_n) = \frac{P(a_1, \dots, a_n | v_i) P(v_i)}{P(a_1, \dots, a_n)}$$

O nosso objetivo é descobrir qual é o v_i que maximiza $P(v_i | a_1, \dots, a_n)$, isto é, pretendemos descobrir

$$\begin{aligned} \arg \max_{v_i \in V} P(v_i | a_1, \dots, a_n) &= \arg \max_{v_i \in V} \frac{P(a_1, \dots, a_n | v_i) P(v_i)}{P(a_1, \dots, a_n)} \\ &= \arg \max_{v_i \in V} P(a_1, \dots, a_n | v_i) P(v_i) \end{aligned}$$

Reparem que não precisamos de dividir por $P(a_1, \dots, a_n)$ visto que o termo é o mesmo para todos os v_i e por isso não é necessário. O problema é que o custo de estimação de $P(a_1, \dots, a_n | v_i)$ é muito elevado, precisaríamos de muitos dados para o conseguir. Assim, este algoritmo assume que os valores dos atributos são condicionalmente independentes, e assim:

$$P(a_1, \dots, a_n | v_i) \approx \prod_{k=1}^n P(a_k | v_i)$$

Temos assim o Classificador Naive Bayes (NB):

$$v_{\text{NB}} = \arg \max_{v_i \in V} P(v_i) \prod_{k=1}^n P(a_k | v_i)$$

Outlook	Temperature	Humidity	Wind	Play Tennis?
Overcast	Hot	Normal	Weak	Yes
Overcast	Mild	High	Strong	Yes
Sunny	Mild	Normal	Strong	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Cool	Normal	Weak	Yes
Overcast	Cool	Normal	Strong	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Mild	High	Weak	Yes
Overcast	Hot	High	Weak	Yes
Rain	Cool	Normal	Strong	No
Sunny	Hot	High	Strong	No
Sunny	Hot	High	Weak	No
Rain	Mild	High	Strong	No
Sunny	Mild	High	Weak	No

Tabela 1: Exemplo de um dataset

Exemplo

O conjunto de dados apresentado na tabela 1 refere-se à decisão de jogar ténis de acordo com vários atributos, Outlook, Temperature, Humidity e Wind. Se estivermos a utilizar um classificador Naive Bayes qual é a classe que ele sugere para o caso dado a seguir:

Outlook	Temperature	Humidity	Wind
Sunny	Cool	High	Strong

Neste caso, o que pretendemos é descobrir v_{NB} de acordo com a formula dada acima:

$$v_{NB} = \arg \max_{v_i \in \{Yes, No\}} P(v_i) \times P(Outlook = Sunny | v_i) \times P(Temperature = Cool | v_i) \times P(Humidity = High | v_i) \times P(Wind = Strong | v_i)$$

Para isso vamos calcular as frequências dos valores de cada atributo para cada um dos v_i tal como se vê na tabela 2. A seguir calculamos as frequências relativas (tabela 3) e podemos finalmente calcular os valores esperados:

Outlook	Yes	No	Temperature	Yes	No	Humidity	Yes	No	Wind	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	Strong	3	3
Overcast	4	0	Mild	4	2	Normal	6	1	Weak	6	2
Rain	3	2	Cool	3	1						
Total	9	5	Total	9	5	Total	9	5	Total	9	5

Tabela 2: Frequências absolutas do dataset da tabela 1

Outlook	Yes	No	Temperature	Yes	No	Humidity	Yes	No	Wind	Yes	No
Sunny	$\frac{2}{9}$	$\frac{3}{5}$	Hot	$\frac{2}{9}$	$\frac{2}{5}$	High	$\frac{3}{9}$	$\frac{4}{5}$	Strong	$\frac{3}{9}$	$\frac{3}{5}$
Overcast	$\frac{4}{9}$	$\frac{0}{5}$	Mild	$\frac{4}{9}$	$\frac{2}{5}$	Normal	$\frac{6}{9}$	$\frac{1}{5}$	Weak	$\frac{6}{9}$	$\frac{2}{5}$
Rain	$\frac{3}{9}$	$\frac{2}{5}$	Cool	$\frac{3}{9}$	$\frac{1}{5}$						

Tabela 3: Frequências relativas do dataset da tabela 1

$$\begin{aligned}
P(\text{PlayTennis} = \text{Yes}) &\times P(\text{Outlook} = \text{Sunny} \mid \text{PlayTennis} = \text{Yes}) \times \\
&P(\text{Temperature} = \text{Cool} \mid \text{PlayTennis} = \text{Yes}) \times \\
&P(\text{Humidity} = \text{High} \mid \text{PlayTennis} = \text{Yes}) \times \\
&P(\text{Wind} = \text{Strong} \mid \text{PlayTennis} = \text{Yes}) = \\
&= \frac{9}{14} \times \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} = 0.00529
\end{aligned}$$

$$\begin{aligned}
P(\text{PlayTennis} = \text{No}) &\times P(\text{Outlook} = \text{Sunny} \mid \text{PlayTennis} = \text{No}) \times \\
&P(\text{Temperature} = \text{Cool} \mid \text{PlayTennis} = \text{No}) \times \\
&P(\text{Humidity} = \text{High} \mid \text{PlayTennis} = \text{No}) \times \\
&P(\text{Wind} = \text{Strong} \mid \text{PlayTennis} = \text{No}) = \\
&= \frac{5}{14} \times \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} = 0.02057
\end{aligned}$$

Sendo assim, neste caso decidiríamos não ir jogar ténis. Para além disso, podemos estimar a probabilidade condicional do valor ser *não* dados os atributos observados:

$$\frac{0.02057}{0.02057 + 0.00529} = 79.54\%$$

Estimar probabilidades

Até agora estimámos as probabilidades pelas frequências relativas simples. Isto é, para estimar a probabilidade $P(\text{Humidity} = \text{Normal} \mid \text{PlayTennis} = \text{No})$ dividimos o número de ocorrências de $\text{Humidity} = \text{Normal}$ quando $\text{PlayTennis} = \text{No}$ pelo número de ocorrências de $\text{PlayTennis} = \text{No}$ o que nos dá $P(\text{Humidity} = \text{Normal} \mid \text{PlayTennis} = \text{No}) = \frac{1}{5}$.

O problema desta estimação é quando a probabilidade real é muito baixa, por exemplo 0.02. Neste caso, é muito provável que a frequência relativa seja 0, e ao multiplicar 0 por outros valores

obtém-se sempre 0. Nestes casos, costuma-se utilizar a **estimativa m da probabilidade**:

$$\frac{n_c + m \times p}{n + m}$$

em que n_c é o número de ocorrências que pretendemos (no nosso caso, quando $Humidity = Normal$ se $PlayTennis = No$) e n é o número total de exemplos para os quais $PlayTennis = No$, p é a *estimativa à priori da probabilidade* que pretendemos determinar e m é uma constante chamada *tamanho equivalente da amostra* que determina o nosso grau de confiança em p relativamente à nossa amostra. Tipicamente escolhemos o valor de p assumindo a distribuição uniforme. Assim, caso existam k valores possíveis, assumimos que $p = \frac{1}{k}$. Por exemplo, ao estimar o valor de $P(Humidity = Normal | PlayTennis = No)$ sabemos que $Humidity$ tem dois valores possíveis (*High* e *Normal*) e por isso $p = \frac{1}{2}$. Se m for zero então a estimativa m da probabilidade dá simplesmente $\frac{n_c}{n}$. Se m for 2 e assumindo para p o valor à priori segundo a distribuição uniforme, a estimativa m daria

$$\frac{n_c + 1}{n + 2}$$

A razão porque se chama a m o tamanho equivalente da amostra é porque aos n valores observados na amostra se juntam m valores virtuais segundo a probabilidade p .

Aprendendo a classificar texto

O NB é um método bastante bem sucedido para classificar textos. Neste caso, o que se faz mais frequentemente é contar as frequências de cada palavra p_k pertencente ao texto que se pretende classificar e calcular a probabilidade de pertencer à classe c_i mediante a fórmula:

$$P(c_i | p_1, \dots, p_n) = P(c_i) \times \prod_{k=1}^n P(p_k | c_i)$$

em que a probabilidade de pertencer à classe c_i é dada pela frequência relativa dos n_i textos pertencentes à classe c_i sobre todos os n textos usados no treino

$$P(c_i) = \frac{n_i}{n}$$

e a probabilidade da palavra p_k aparecer nos textos da classe c_i é dada pela estimativa m em que m é o número de palavras diferentes que ocorrem nos textos que foram usados no treino sendo dada por

$$P(p_k | c_i) = \frac{n_k + 1}{n + |Vocabulario|}$$

em que n_k é o número total de vezes em que a palavra p_k aparece nos textos da classe c_i , n é o número total de palavras que ocorre nos textos da classe c_i e $|Vocabulario|$ é o número total de palavras diferentes que ocorre nos textos.

Agora aplicamos o algoritmo para descobrir a classe:

$$c_{NB} = \arg \max_{i \in classes} P(c_i) \times \prod_{k=1}^n P(p_k | c_i)$$

Para evitar as probabilidades muito baixas, que podem tornar a multiplicação dos valores zero e por isso impedir o algoritmo de funcionar corretamente, utilizam-se logaritmos:

$$c_{NB} = \arg \max_{i \in \text{classes}} \log P(c_i) + \sum_{k=1}^n \log P(p_k | c_i)$$

Vamos ver um exemplo para nos ajudar a perceber.

texto	classe
a baixa do porto	Porto
o mercado do bolhão é no porto	Porto
a câmara do porto fica no centro do porto	Porto
a baixa de lisboa	Lisboa
o porto de lisboa	Lisboa

Vocabulário 14

	Porto	Lisboa
Exemplos	3	2
Palavras	20	8
baixa	1	1
porto	4	1
mercado	1	0
bolhão	1	0
câmara	1	0
baixa	1	1
lisboa	0	2

$$P(\text{Porto}) = \frac{3}{5}$$

$$P(\text{Lisboa}) = \frac{2}{5}$$

$$P(\text{porto}|\text{Porto}) = \frac{4+1}{20+14} = 0.1471$$

$$P(\text{porto}|\text{Lisboa}) = \frac{1+1}{8+14} = 0.0909$$

$$p(\text{mercado}|\text{Porto}) = \frac{1+1}{20+14} = 0.0588$$

$$p(\text{mercado}|\text{Lisboa}) = \frac{0+1}{8+14} = 0.0455$$

E assim, ao tentar classificar o texto *porto porto porto mercado* teríamos:

$$\begin{aligned} \text{Porto} \quad & P(\text{Porto}) \times p(\text{porto}|\text{Porto})^3 \times P(\text{mercado}|\text{Porto}) = 0.6 \times 0.1471^3 \times 0.0588 = 0.0001123 \\ \text{Lisboa} \quad & P(\text{Lisboa}) \times p(\text{porto}|\text{Lisboa})^3 \times P(\text{mercado}|\text{Lisboa}) = 0.4 \times 0.0909^3 \times 0.0455 = 0.0000137 \end{aligned}$$

O que nos permitiria concluir que estavamos a falar sobre o Porto com probabilidade

$$\frac{0.0001123}{0.0001123 + 0.0000137} = 89.13\%$$

Técnicas específicas para deteção de spam

Quando se pretende detetar spam calcula-se a probabilidade de uma dada palavra P pertencer a um texto de spam da seguinte forma:

$$P(S|P) = \frac{P(P|S) \times P(S)}{P(P|S) \times P(S) + P(P|N) \times P(N)}$$

em que S corresponde a spam, N a não spam, e P a uma dada palavra. Caso se assuma que $P(S) = P(N) = 0.5$ (na verdade as estimativas indicam que $P(S) \approx 0.8$), a fórmula pode ser simplificada da seguinte forma:

$$P(S|P) = \frac{P(P|S)}{P(P|S) + P(P|N)}$$

A probabilidade de um texto que contém as palavras P_1, \dots, P_n ser spam é dada por:

$$P(S|P_1, \dots, P_n) = \frac{\prod_{i=1}^n P(S|P_i)}{\prod_{i=1}^n P(S|P_i) + \prod_{i=1}^n P(N|P_i)}$$

Repare que como só existem as classes S e N , a probabilidade $P(N|P_i) = 1 - P(S|P_i)$.