

Análise de Dados e Inferência Estatística com Python



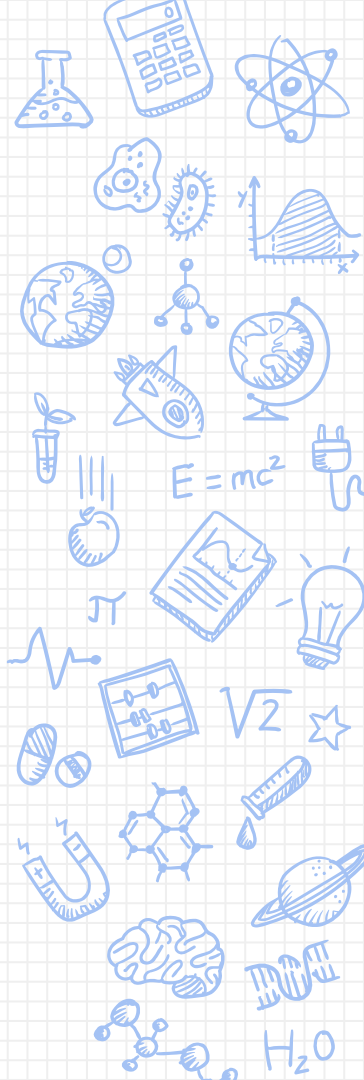
-

Introdução

É uma técnica de classificação baseada no teorema de Bayes com uma suposição de independência entre os preditores.

Em termos simples, um classificador Naive Bayes assume que a presença de uma característica particular em uma classe não está relacionada com a presença de qualquer outro recurso.

Por exemplo, um fruto pode ser considerado como uma maçã se é vermelho, redondo, e tiver cerca de 3 polegadas de diâmetro. Mesmo que esses recursos dependam uns dos outros ou da existência de outras características, todas estas propriedades contribuem de forma independente para a probabilidade de que este fruto é uma maçã e é por isso que é conhecido como 'Naive' (ingênuo).



O Teorema de Bayes fornece uma forma de calcular a probabilidade posterior $P(C|X)$ a partir de $P(C)$, $P(x)$ e $P(X|c)$. Veja a equação abaixo:

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Naive Bayes

- $P(c | x)$ é a probabilidade posterior da classe $\{c, \text{alvo}\}$ dada preditor $\{x, \text{atributos}\}$.
- $P(c)$ é a probabilidade original da classe.
- $P(x | c)$ é a probabilidade que representa a probabilidade de preditor dada a classe.
- $P(x)$ é a probabilidade original do preditor.

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

Diagram illustrating the components of the Naive Bayes formula:

- $P(c | x)$ is labeled "Probabilidade posterior" (Posterior Probability).
- $P(x | c)$ is labeled "Preditor da probabilidade posterior" (Predictor of the posterior probability).
- $P(c)$ is labeled "Probabilidade original da Classe" (Original Probability of the Class).
- $P(x)$ is labeled "Probabilidade" (Probability).

$$P(c | X) = P(x_1 | c) \times P(x_2 | c) \times \dots \times P(x_n | c) \times P(c)$$



Como o algoritmo Naive Bayes funciona?

TEMPO	"PLAY"
Sol	Não
Nublado	Sim
Chuva	Sim
Sol	Sim
Sol	Sim
Nublado	Sim
Chuva	Não
Chuva	Não
Sol	Sim
Chuva	Sim
Sol	Não
Nublado	Sim
Nublado	Sim
Chuva	Não

Tabela de frequência		
Clima	Não	Sim
Nublado	0	4
Sol	3	2
Chuva	2	3
Total	5	9

Tabela de probabilidade		
Clima	Não	Sim
Nublado	0	4
Sol	3	2
Chuva	2	3
Total	5	9
	$\approx 5/14$	$\approx 9/14$
	0,36	0,64

$\approx 4/14$	0,29
$\approx 5/14$	0,36
$\approx 5/14$	0,36

Problema: Os jogadores irão jogar se o tempo está ensolarado. Esta afirmação está correta?

Como o algoritmo Naive Bayes funciona?

Podemos resolver isso usando o método discutido acima de probabilidade posterior.

$$P(\text{Sim} | \text{Ensolarado}) = P(\text{Ensolarado} | \text{Sim}) * P(\text{Sim}) / P(\text{Ensolarado})$$

Aqui temos $P(\text{Ensolarado} | \text{Sim}) = 3/9 = 0,33$, $P(\text{Ensolarado}) = 5/14 = 0,36$, $P(\text{Sim}) = 9/14 = 0,64$

Agora, $P(\text{Sim} | \text{Ensolarado}) = 0,33 * 0,64 / 0,36 = 0,60$, que tem maior probabilidade.

Naive Bayes usa um método similar para prever a probabilidade de classe diferente com base em vários atributos. Este algoritmo é usado principalmente em classificação de texto e com os problemas que têm múltiplas classes.



Quais são os prós e contras de Naive Bayes?

Prós:

- É fácil e rápido para prever o conjunto de dados da classe de teste. Também tem um bom desempenho na previsão de classes múltiplas.
- Quando a suposição de independência prevalece, um classificador Naive Bayes tem melhor desempenho em comparação com outros modelos como regressão logística, e você precisa de menos dados de treinamento.
- O desempenho é bom em caso de variáveis categóricas de entrada comparada com variáveis numéricas. Para variáveis numéricas, assume-se a distribuição normal (curva de sino, que é uma suposição forte).



Contras:

- Se a variável categórica tem uma categoria (no conjunto de dados de teste) que não foi observada no conjunto de dados de treinamento, então o modelo irá atribuir uma probabilidade de 0 (zero) e não será capaz de fazer uma previsão. Isso é muitas vezes conhecido como “Zero Frequency”. Para resolver isso, podemos usar a técnica de alisamento. Uma das técnicas mais simples de alisamento é a chamada estimativa de Laplace.
- Por outro lado, Naive Bayes é também conhecido como um mau estimador, por isso, as probabilidades calculadas não devem ser levadas muito a sério.
- Outra limitação do Naive Bayes é a suposição de preditores independentes. Na vida real, é quase impossível ter um conjunto de indicadores que sejam completamente independentes.

-

Gaussian: É usado na classificação e assume uma distribuição normal.

Com base no seu conjunto de dados, você pode escolher qualquer um dos modelos acima discutidos.

Aqui vão algumas dicas para melhorar o poder do Modelo Naive Bayes:

- Se as funções contínuas não têm distribuição normal, devemos usar a transformação ou métodos diferentes para convertê-las na distribuição normal.
- Se o conjunto de dados de teste tem problema de frequência zero, aplique a técnica de suavização “Laplace Correction” para prever a classe de conjunto de dados de teste.
- Remova características correlacionadas, já que as características altamente correlacionadas são votadas duas vezes no modelo e podem levar a um excesso de importância.
- Classificadores Naive Bayes têm opções limitadas para ajuste de parâmetros como $\alpha = 1$ para alisamento, `fit_prior = [Verdade | Falso]` para saber a classe de probabilidades anteriores ou não e algumas outras opções. Recomenda-se focar no pré-processamento de dados e seleção de recursos.
- Você poderia pensar em aplicar alguma técnica de combinação de classificador como “ensembling”, “bagging” e “boosting”, mas na prática esses métodos não ajudariam, pois sua finalidade é reduzir a variância. Naive Bayes não tem variância para minimizar.



Obrigado!