



UNIVERSIDADE FEDERAL DO RIO GRANDE DO NORTE
CENTRO DE TECNOLOGIA
DEPARTAMENTO DE ENGENHARIA ELÉTRICA
ELE0606 - TÓPICOS ESPECIAIS EM INTELIGÊNCIA ARTIFICIAL

Docente:

José Alfredo Ferreira Costa

Autor:

Vinícius Venceslau Venancio da Penha

Métodos de Classificação

Natal - RN
15 de outubro de 2023

RESUMO

Este trabalho explora a Regra de Bayes, uma técnica estatística fundamental com aplicações em estatística, aprendizado de máquina e inteligência artificial. A Regra de Bayes permite o cálculo de probabilidades condicionais, atualizando crenças com base em novas evidências. O documento aborda os princípios da Regra de Bayes e fornece um exemplo intuitivo.

Ademais, descreve a aplicação da Regra de Bayes no Classificador Bayesiano, com foco em Estimação Paramétrica e Estimação Não-Paramétrica. O trabalho também introduz o Naive Bayes, um método de classificação baseado na suposição de independência condicional entre variáveis preditoras, e apresenta exemplos de sua aplicação em duas bases de dados: "Wine" e "Heart Disease."

Infere-se, portanto, que o trabalho fornece códigos em Python, detalha a organização e documentação adequada do código, e conclui com uma visão geral dos tipos comuns de Naive Bayes.

Sumário

1	Regra de Bayes	2
1.1	Princípios	2
1.2	Simulação	4
1.2.1	Base de Dados - Wine	5
1.2.2	Base de Dados - Heart Disease	7
2	Regra Simplificada - Naive Bayes	9
2.1	Princípios	9
2.2	Simulação	9
2.2.1	Base de Dados - Wine	10
2.2.2	Base de Dados - Heart Disease	11
3	Conclusão	13
4	Referências	14

1 Regra de Bayes

A Lei de Bayes, em homenagem ao matemático Thomas Bayes, é uma técnica estatística fundamental com amplas aplicações em áreas como estatística, aprendizado de máquina e inteligência artificial.

Esse método fornece uma poderosa maneira de calcular probabilidades condicionais, permitindo que crenças sejam atualizadas com base em novas evidências. Nesse viés, o teorema desempenha um papel central em muitos campos, incluindo classificação, diagnóstico médico, filtragem de spam e reconhecimento de padrões.

1.1 Princípios

A Regra de Bayes é fundamentada em princípios estatísticos sólidos e é expressa na seguinte equação:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \quad (1)$$

Nesse sentido, têm-se:

- $P(A|B)$ é a probabilidade posterior do evento A ocorrer, dadas as evidências B.
- $P(B|A)$ é a probabilidade da evidência B ocorrer, supondo que o evento A seja verdadeiro.
- $P(A)$ é a probabilidade prévia do evento A, antes de considerar as evidências B.
- $P(B)$ é a probabilidade da evidência B ocorrer, independentemente do evento A.

Diante disso, a seguir, será exibido um exemplo intuitivo a respeito da metodologia descrita anteriormente.

Exemplo Intuitivo:

Suponha a existência de um teste médico para uma doença rara, e deseja-se calcular a probabilidade de uma pessoa realmente ter a doença, dado que o teste foi positivo.

Nesse cenário, $P(A|B)$ representaria a probabilidade da pessoa ter a doença após um resultado positivo, $P(B|A)$ seria a probabilidade do teste ser positivo se a pessoa tiver a doença, $P(A)$ seria a probabilidade prévia da pessoa ter a doença e $P(B)$ seria a probabilidade de um resultado positivo, independentemente da presença da doença.

Além disso, com o intuito de produzir uma melhor visualização e/ou simplificação matemática do método em análise, pode-se expressar $P(A|B)$ da seguinte maneira:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Uma vez que $P(B|A) \cdot P(A)$ é igual a interseção entre os conjuntos A e B, ou seja, $P(A \cap B)$.

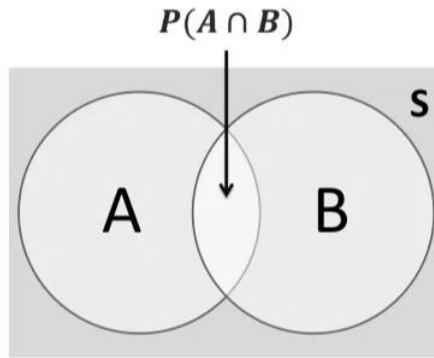


Figura 1: Interseção entre dois conjuntos.

Adicionalmente, o **Classificador Bayesiano** é uma técnica de classificação de dados amplamente utilizada no campo da ciência de dados. Ele se baseia na Teoria da Decisão Bayesiana, que é uma abordagem estatística para a classificação de padrões.

Diante disso, a principal ideia por trás do Classificador Bayesiano é usar a Regra de Bayes para calcular a probabilidade de um determinado ponto de dados pertencer a uma classe específica com base nas informações disponíveis.

Em suma, a abordagem Bayesiana considera a probabilidade condicional das classes, dadas as observações (dados). Ele calcula a probabilidade de que uma determinada instância de dados pertença a uma classe específica, considerando as probabilidades prévias das classes e as probabilidades das observações condicionadas a cada classe.

Estimação Paramétrica:

A classificação Bayesiana pode ser realizada de duas maneiras, e uma delas é a **Estimação Paramétrica**. Neste método, avalia-se o formato da distribuição dos dados em observação. O modelo de disposição dos dados pode seguir distribuições como Normal, Weibull, Exponencial, entre outras.

Portanto, na Estimação Paramétrica, o primeiro passo é determinar o padrão de distribuição mais adequado para os dados. Isso é feito através da análise exploratória dos dados e testes estatísticos. Uma vez que o modelo de distribuição é selecionado, o próximo passo é estimar os parâmetros dessa distribuição. Os parâmetros são valores que definem a forma da distribuição, como a média e o desvio padrão em uma distribuição normal.

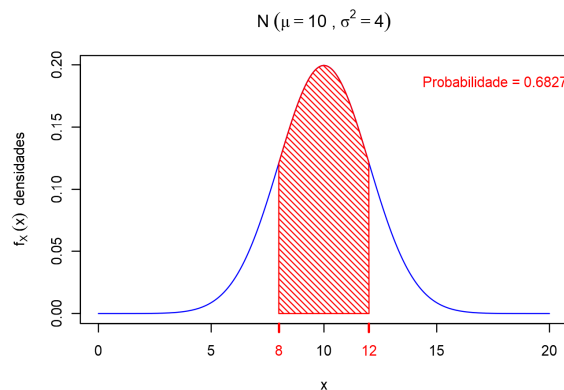


Figura 2: Exemplo de distribuição Normal.

Estimação Não-Paramétrica:

A outra forma de classificação é a **Estimação Não-Paramétrica**. Neste caso, o Classificador Bayesiano estima as probabilidades diretamente dos dados de treinamento, sem assumir um modelo de distribuição específico.

Esse método é útil quando os dados são complexos e não podem ser facilmente modelados por uma distribuição paramétrica. Outrossim, a Estimação Não-Paramétrica é mais flexível, pois não restringe a forma da distribuição dos dados, permitindo que o classificador se adapte a uma variedade de situações.

Por fim, ambos os métodos têm suas vantagens e desvantagens, e a escolha entre eles depende da natureza dos dados e das suposições que pode-se fazer sobre a distribuição subjacente.

Por conseguinte, em estatística e análise de dados, o termo "kernel" refere-se a uma função utilizada na estimação não paramétrica de densidade de probabilidade ou na suavização de dados.

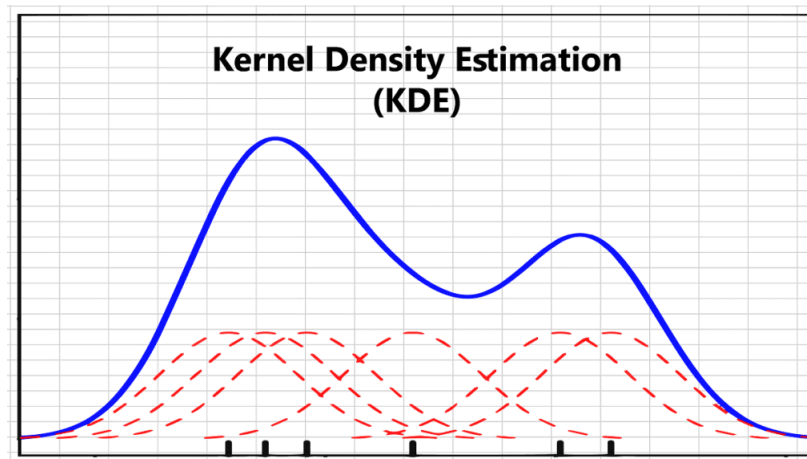


Figura 3: Ilustração da função Kernel.

1.2 Simulação

Nas subseções abaixo, serão abordadas as implementações práticas dos princípios previamente discutidos neste relatório. Por meio da utilização de código Python, será ilustrada a aplicação efetiva desses conceitos em cenários do mundo real. A simulação proporcionará uma exposição concreta das teorias e estratégias discutidas, demonstrando de forma palpável a sua relevância e aplicabilidade.

Nessa conjuntura, vale salientar que todos os códigos desenvolvidos pelo discente, encontram-se disponibilizados na plataforma do Github no repositório adequado, uma vez que se tem o intuito de facilitar a visualização dos programas produzidos.

Por fim, todo material de apoio usado para confeccionar os códigos estão expostos na referência deste documento.

1.2.1 Base de Dados - Wine

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.datasets import load_wine
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6 from scipy.stats import multivariate_normal
7 from sklearn.metrics import accuracy_score, classification_report,
  confusion_matrix
8
9 #Carregar o conjunto de dados Wine:
10 wine_data = load_wine()
11 wine_df = pd.DataFrame(wine_data['data'], columns=wine_data['
  feature_names'])
12 wine_df['classe'] = wine_data['target'] # Adicionar a coluna de classe
13
14 #Preparar os dados:
15 X = wine_df.drop('classe', axis=1).to_numpy()
16 y = wine_df['classe'].to_numpy()
17
18 #Normalizar os dados:
19 scaler = StandardScaler()
20 X = scaler.fit_transform(X)
21
22 #Dividir os dados em conjuntos de treinamento e teste:
23 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
  random_state=42)
24
25 #Inicializar a matriz de probabilidades:
26 num_classes = len(np.unique(y))
27 probabilidadeDaClasse = np.zeros(num_classes)
28
29 for i in range(num_classes):
30     indiceDaClasse = np.where(y_train == i)
31     probabilidadeDaClasse[i] = len(indiceDaClasse[0]) / len(y_train)
32
33 #Matriz para armazenar as probabilidades de classe para cada observacao
  de teste:
34 matrizProbabilidade = np.zeros((X_test.shape[0], num_classes))
35
36 #Calcular as probabilidades de classe para as observacoes de teste:
37 for i in range(num_classes):
38     indiceDaClasse = np.where(y_train == i)
39     dadosDeTreino = X_train[indiceDaClasse]
40     media = np.mean(dadosDeTreino, axis=0)
41     covariancia = np.cov(dadosDeTreino, rowvar=False)
42
43     for j in range(X_test.shape[0]):
44         observation = X_test[j]
45         matrizProbabilidade[j, i] = multivariate_normal.pdf(observation,
  mean=media, cov=covariancia, allow_singular=True) *
  probabilidadeDaClasse[i]
46
47 #Prever as classes com base nas probabilidades:
48 y_pred = np.argmax(matrizProbabilidade, axis=1)
49
50 #Calcular a acuracia das previsoes:
```

```

51 accuracy = accuracy_score(y_test, y_pred)
52 print('Acuracia:', accuracy)
53
54 #Gerar Matriz Confusao:
55 confusion = confusion_matrix(y_test, y_pred)
56 report = classification_report(y_test, y_pred, target_names=['Classe 0',
    'Classe 1', 'Classe 2'])
57
58 print("\nMatriz de Confusao:\n", confusion)
59 print("\nRelatorio de Classificacao:\n", report)

```

Listing 1: Classificador Bayesiano - Wine.

```

1 Acuracia: 0.9722
2
3 Matriz de Confusao:
4 [[14  0  0]
5  [ 0 14  0]
6  [ 0  1  7]]
7
8 Relatorio de Classificacao:
9
10      precision    recall  f1-score   support
11
12  Classe 0       1.00      1.00      1.00        14
13  Classe 1       0.93      1.00      0.97        14
14  Classe 2       1.00      0.88      0.93         8
15
16  accuracy              0.97        36
17  macro avg       0.98      0.96      0.97        36
18  weighted avg    0.97      0.97      0.97        36

```

Listing 2: Saída 1.

1.2.2 Base de Dados - Heart Disease

```
1 #Permitir o google colab acessar os arquivos do Drive:
2 from google.colab import drive
3 drive.mount('/content/drive')
4
5 import pandas as pd
6 import numpy as np
7 from sklearn.model_selection import train_test_split
8 from sklearn.preprocessing import StandardScaler
9 from scipy.stats import multivariate_normal
10 from sklearn.metrics import accuracy_score, classification_report,
    confusion_matrix
11
12 #Carregamento da base de dados a partir do Google Drive:
13 caminho_arquivo = "/content/drive/My Drive/heart.csv"
14 banco_de_dados = pd.read_csv(caminho_arquivo)
15
16 #Criar um Dataframe:
17 heart_df = banco_de_dados
18
19 #Preparar os dados:
20 X = heart_df.drop('target', axis=1).to_numpy()
21 y = heart_df['target'].to_numpy()
22
23 #Normalizar os dados:
24 scaler = StandardScaler()
25 X = scaler.fit_transform(X)
26
27 #Dividir os dados em conjuntos de treinamento e teste:
28 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
29
30 #Inicializar a matriz de probabilidades:
31 num_classes = len(np.unique(y))
32 probabilidadeDaClasse = np.zeros(num_classes)
33
34 for i in range(num_classes):
35     indiceDaClasse = np.where(y_train == i)
36     probabilidadeDaClasse[i] = len(indiceDaClasse[0]) / len(y_train)
37
38 #Matriz para armazenar as probabilidades de classe para cada observacao
    de teste:
39 matrizProbabilidade = np.zeros((X_test.shape[0], num_classes))
40
41 #Calcular as probabilidades de classe para as observacoes de teste:
42 for i in range(num_classes):
43     indiceDaClasse = np.where(y_train == i)
44     dadosDeTreino = X_train[indiceDaClasse]
45     media = np.mean(dadosDeTreino, axis=0)
46     covariancia = np.cov(dadosDeTreino, rowvar=False)
47
48     for j in range(X_test.shape[0]):
49         observation = X_test[j]
50         matrizProbabilidade[j, i] = multivariate_normal.pdf(observation,
            mean=media, cov=covariancia, allow_singular=True) *
            probabilidadeDaClasse[i]
51
52 #Prever as classes com base nas probabilidades:
```

```

53 y_pred = np.argmax(matrizProbabilidade, axis=1)
54
55 #Calcular a acuracia das previsoes:
56 accuracy = accuracy_score(y_test, y_pred)
57 print(f'Acuracia: {accuracy:.4f}')
58
59 #Gerar Matriz Confusao:
60 confusion = confusion_matrix(y_test, y_pred)
61 report = classification_report(y_test, y_pred, target_names=['Classe 0',
        'Classe 1'])
62
63 print("\nMatriz de Confusao:\n", confusion)
64 print("\nRelatorio de Classificacao:\n", report)

```

Listing 3: Classificador Bayesiano - Heart Disease.

```

1 Drive already mounted at /content/drive; to attempt to forcibly remount,
  call drive.mount("/content/drive", force_remount=True).
2 Acuracia: 0.8634
3
4 Matriz de Confusao:
5 [[85 21]
6  [ 7 92]]
7
8 Relatorio de Classificacao:
9
10      precision    recall  f1-score   support
11
12  Classe 0       0.92      0.80      0.86       106
13  Classe 1       0.81      0.93      0.87        99
14
15   accuracy              0.86       205
16  macro avg       0.87      0.87      0.86       205
17  weighted avg       0.87      0.86      0.86       205

```

Listing 4: Saída 2.

2 Regra Simplificada - Naive Bayes

2.1 Princípios

Naive Bayes é frequentemente usado para tarefas de classificação, onde o objetivo é prever a classe de um objeto ou evento com base em informações disponíveis. A fórmula geral para classificação Bayesiana é:

$$P(C_k|X) = \frac{P(X|C_k) \cdot P(C_k)}{P(X)} \quad (2)$$

Onde:

- $P(C_k|X)$ é a probabilidade da classe C_k dado um conjunto de características X .
- $P(X|C_k)$ é a probabilidade de observar X dado que a classe é C_k .
- $P(C_k)$ é a probabilidade prévia da classe C_k .
- $P(X)$ é a probabilidade de observar X .

O termo "Naive" em Naive Bayes se refere à suposição de independência condicional entre as variáveis preditoras (características) dadas as classes.

Dessa forma, isso significa que a presença ou valor de uma característica não depende das outras características quando a classe é conhecida. Essa suposição simplifica o cálculo das probabilidades, tornando-o mais viável computacionalmente.

Existem vários tipos de Naive Bayes, sendo os mais comuns:

1. **Naive Bayes Gaussiano:** Assume que as características seguem uma distribuição normal (Gaussiana).
2. **Naive Bayes Multinomial:** Usado para dados discretos, como contagens de palavras em um documento. É amplamente usado em classificação de texto.
3. **Naive Bayes Bernoulli:** Usado para características binárias, onde cada característica é representada como uma variável booleana.

2.2 Simulação

Os códigos em Python estão sendo apresentados, com ênfase na organização da formatação, utilização de comentários esclarecedores e adoção de nomes descritivos para variáveis, juntamente com a estruturação do código em blocos lógicos por meio de funções e classes para facilitar a manutenção e compreensão.

Assim, a documentação adequada desempenha um papel importante na promoção da colaboração eficaz e na comunicação clara de soluções e resultados em projetos Python.

2.2.1 Base de Dados - Wine

```
1 from sklearn.datasets import load_wine
2 from sklearn.model_selection import train_test_split
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.metrics import accuracy_score, classification_report,
  confusion_matrix
5
6 wine = load_wine()
7 X = wine.data
8 y = wine.target
9
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4,
  random_state=13)
11
12 nb_classifier = GaussianNB()
13 nb_classifier.fit(X_train, y_train)
14
15 y_pred = nb_classifier.predict(X_test)
16
17 accuracy = accuracy_score(y_test, y_pred)
18 confusion = confusion_matrix(y_test, y_pred)
19 report = classification_report(y_test, y_pred, target_names=wine.
  target_names)
20
21 print("Acuracia:", accuracy)
22 print("\nMatriz de Confusao:\n", confusion)
23 print("\nRelatorio de Classificacao:\n", report)
```

Listing 5: Código Naive Bayes - Wine.

```
1 Acuracia: 0.9861111111111112
2
3 Matriz de Confusao:
4 [[24  0  0]
5  [ 0 30  1]
6  [ 0  0 17]]
7
8 Relatorio de Classificacao:
9               precision    recall  f1-score   support
10
11    class_0             1.00      1.00      1.00        24
12    class_1             1.00      0.97      0.98        31
13    class_2             0.94      1.00      0.97        17
14
15    accuracy                   0.99        72
16    macro avg              0.98      0.99      0.99        72
17    weighted avg           0.99      0.99      0.99        72
```

Listing 6: Saída 3.

2.2.2 Base de Dados - Heart Disease

```
1 #Permitir o google colab acessar os arquivos do Drive:
2 from google.colab import drive
3 drive.mount('/content/drive')
4
5 #Importar bibliotecas:
6 import numpy as np
7 import pandas as pd
8 import matplotlib.pyplot as plt
9 from sklearn.model_selection import train_test_split
10 from sklearn.metrics import accuracy_score, classification_report,
    confusion_matrix
11
12 #Carregamento da base de dados a partir do Google Drive:
13 caminho_arquivo = "/content/drive/My Drive/heart.csv"
14 banco_de_dados = pd.read_csv(caminho_arquivo)
15
16 #Criar um Dataframe:
17 heart_df = banco_de_dados
18
19 #Armazenamento dos rotulos de classe em uma variavel:
20 heart_classe = heart_df['target']
21
22 #Remocao da coluna 'target' do DataFrame:
23 heart_df.drop(['target'], axis=1, inplace=True)
24
25 #Divisao dos dados em conjuntos de treinamento e teste:
26 X_train, X_test, y_train, y_test = train_test_split(heart_df,
    heart_classe, test_size=0.3, random_state=13)
27
28 #Criar um classificador Naive Bayes Gaussiano:
29 nb_classifier = GaussianNB()
30
31 #Treinar o classificador:
32 nb_classifier.fit(X_train, y_train)
33
34 #Fazer previsoes no conjunto de teste:
35 y_pred = nb_classifier.predict(X_test)
36
37 #Avaliar o desempenho do classificador:
38 accuracy = accuracy_score(y_test, y_pred)
39 confusion = confusion_matrix(y_test, y_pred)
40 report = classification_report(y_test, y_pred, target_names=['Classe 0',
    'Classe 1'])
41
42 print("Acuracia:", accuracy)
43 print("\nMatriz de Confusao:\n", confusion)
44 print("\nRelatorio de Classificacao:\n", report)
```

Listing 7: Código Naive Bayes - Heart Disease.

```

1 Drive already mounted at /content/drive; to attempt to forcibly remount,
  call drive.mount("/content/drive", force_remount=True).
2 Acuracia: 0.814935064935065
3
4 Matriz de Confusao:
5 [[113  39]
6  [ 18 138]]
7
8 Relatorio de Classificacao:
9
10      precision      recall  f1-score   support
11
12  Classe 0          0.86         0.74         0.80         152
13  Classe 1          0.78         0.88         0.83         156
14
15  accuracy                    0.81         308
16  macro avg          0.82         0.81         0.81         308
17  weighted avg          0.82         0.81         0.81         308

```

Listing 8: Saída 4.

3 Conclusão

A Regra de Bayes é uma ferramenta valiosa que fornece uma estrutura sólida para calcular probabilidades condicionais em diversas aplicações. Ela é fundamental para entender e implementar algoritmos de aprendizado de máquina, como o Naive Bayes, que se baseiam em princípios bayesianos.

Além disso, a Regra de Bayes é uma das bases da estatística e tem sido amplamente aplicada em diagnóstico médico, previsão de eventos e muito mais. É importante compreender o seu funcionamento e aplicabilidade em diversas situações para utilizar efetivamente esse poderoso método estatístico e de aprendizado de máquina.

4 Referências

LUIZ, E.; OLIVEIRA. **Reconhecimento de Padrões Teoria da Decisão Bayesiana**. [s.l.: s.n.]. Disponível em: <https://www.inf.ufpr.br/lesoliveira/padroes/decisaobayesiana.pdf>. Acesso em: 5 de outubro de 2023.

Ciência de Dados: Teoria da Decisão Bayesiana na classificação de dados. Disponível em: <https://www.youtube.com/watch?v=8zAKWE0dGsg>. Acesso em: 5 de outubro de 2023.

Data Mining Classification: Alternative Techniques. Disponível em: https://www-users.cs.umn.edu/~kumar001/dmbook/slides/chap4_naive_bayes.pptx. Acesso em: 5 de outubro de 2023.

OLIVEIRA, L. **Teorema de Bayes e Probabilidade**. Disponível em: <https://medium.com/data-hackers/teorema-de-bayes-probabilidade-d5ead2df1379>. Acesso em: 14 de outubro de 2023.