

Aprendizado de Máquina - Atividade Prática Naïve Bayes

Willian Reichert (134090)

Análise do questionário

a) Verdadeiro, visto que essa é a ideia do algoritmo: a probabilidade condicional da ocorrência do valor de um atributo dada a ocorrência de uma classe específica leva em consideração apenas a quantidade de vezes que esse valor aparece dentre os dados pertencentes a essa classe.

b) Verdadeiro. Como $P(\text{target} = \text{acc}) = 0.4$ e $P(\text{target} = \text{unacc}) = 0.6$, a probabilidade de uma instância ser da classe *unacc* é maior. Esses valores podem ser calculados rapidamente pela divisão da quantidade de instâncias da cada classe pela quantidade total de instâncias.

c) Verdadeiro. O cálculo é feito pela divisão da quantidade de instâncias de *acc* que possuem o valor *med* no atributo *price* pela quantidade de instâncias de *acc*, o que nos dá 3/10.

d) Parcialmente falso. A instância será classificada pelo modelo como *acc*, porém as probabilidades a posteriori de ambas as classes são diferentes: $P(\text{target} = \text{acc}|x) = 0.0289$ e $P(\text{target} = \text{unacc}|x) = 0.0249$. A implementação do algoritmo para esse dataset específico foi realizada em Python, e o código está no final desse relatório.

e) Verdadeiro. Como não foi utilizada a correção de Laplace, a probabilidade $P(\text{safety} = \text{low}|\text{acc})$ é igual a zero, já que não há nenhuma instância de *acc* com o atributo *safety* igual a *low*.

Código

O código abaixo, implementado para a resolução dos três itens que envolvem cálculos, foi desenvolvido em Python 3.10 com a utilização da biblioteca Pandas. Ele também foi enviado pelo Moodle em um arquivo compactado juntamente com esse relatório.

```
1 from pandas import read_csv, DataFrame
2
3
4 def main() -> None:
5     # leitura dos dados da base para um dataframe
6     data: DataFrame = read_csv('data.csv')
7
8     # item b do questionario
9     p_acc: float = data['target'].value_counts()['acc'] / data.shape[0] # P(target
10    = acc)
11    p_unacc: float = data['target'].value_counts()['unacc'] / data.shape[0] # P(
12    target = unacc)
13    print(f'P(target = acc) = {p_acc}')
14    print(f'P(target = unacc) = {p_unacc}')
15    print(f'P(target = unacc) > P(target = acc)? {p_unacc > p_acc}', end='\n\n')
16
17    # item c do questionario
18    acc_data: DataFrame = data.query('target == "acc"') # filtra os dados da
19    classe acc
20    pcond_price_med: float = acc_data['price'].value_counts()['med'] / acc_data.
21    shape[0] # P(price = med | target = acc)
22    print(f'P(price = med | target = acc) = {pcond_price_med}')
23    print(f'P(price = med | target = acc) == 0.3? {pcond_price_med == 0.3}', end='\n\n')
24
25    # item d do questionario
26    acc_prod: float = 1.0 # produtorio de P(target = acc | x)
27    unacc_prod: float = 1.0 # produtorio de P(target = unacc | x)
28    unacc_data: DataFrame = data.query('target == "unacc"') # filtra os dados da
29    classe unacc
30    test_instance: dict = {'price': 'low', 'lug_boot': 'small', 'safety': 'high'}
31    # instancia de teste x
32    for attr, value in test_instance.items():
33        acc_prod *= acc_data[attr].value_counts()[value] / acc_data.shape[0] #
34        prod * P(attr = value | target = acc)
```

```

28     unacc_prod *= unacc_data[attr].value_counts()[value] / unacc_data.shape[0]
    # prod * P(attr = value | target = unacc)
29     ppost_acc: float = p_acc * acc_prod # P(target = acc | x)
30     ppost_unacc: float = p_unacc * unacc_prod # P(target = unacc | x)
31     print(f'P(target = acc | x) = {ppost_acc:.5f}')
32     print(f'P(target = unacc | x) = {ppost_unacc:.5f}')
33     print(f'P(target = acc | x) > P(target = unacc | x)? {ppost_acc > ppost_unacc}'
    )
34     return
35
36
37 if __name__ == '__main__':
38     main()

```