

Nome: LuãMoreiraPonciano

Pontuação Total da Avaliação: 6.02 pontos

=====

Correção da Questão 1:

Pergunta: 1a) Uma empresa está desenvolvendo um sistema para classificar mensagens recebidas como "Urgente" ou "Não Urgente" com base nas palavras presentes na mensagem. Foi analisado um conjunto de 100 mensagens, e os dados a seguir foram coletados: Mensagens Urgentes: 30; Mensagens Não Urgentes: 70; Palavra Presente "imediatamente": 15 (Mensagens Urgentes) e 5 (Mensagens Não Urgentes); Palavra Presente "problema": 10 (Mensagens Urgentes) e 10 (Mensagens Não Urgentes); Palavra Presente "atraso": 8 (Mensagens Urgentes) e 12 (Mensagens Não Urgentes). Calcule as probabilidades de uma mensagem ser "Urgente" e "Não Urgente" com base no conjunto de dados fornecido. Adicionalmente, determine as probabilidades condicionais para cada palavra ("imediatamente", "problema" e "atraso") em relação às mensagens "Urgentes" e "Não Urgentes".  
Diretriz(es):  $\{('P("urgente")=30/100=0.30', 0.32), ('P("atraso"|"não urgente")=12/70=0.17', 0.32), ('P("imediatamente"|"urgente")=15/30=0.5', 0.32), ('P("imediatamente"|"não urgente")=5/70=0.07', 0.32), ('P("atraso"|"urgente")=8/30=0.27', 0.32), ('P("problema"|"urgente")=10/30=0.33', 0.32), ('P("não urgente")=70/100=0.70', 0.32), ('P("problema"|"não urgente")=10/70=0.14', 0.32)\}$

=====Resposta do estudante da questão 1:

A1) De acordo com os dados apresentados, uma mensagem qualquer tem uma chance de 0.3 (30%) de ser urgente. Em outras palavras,  $U = 30/100$  e  $NU = 70/100$ .

Além disso, para as palavras apresentadas:

A chance da palavra imediatamente aparecer é  $20/100$ , ou seja 0.2 (20%). No total, das 30 mensagens urgentes, 15 continham **imediatamente**, ou seja  $UI = 15/30$  ou 0.5 (50%). Em contrapartida das 70 mensagens não urgentes, apenas 5 continham essa palavra, ou seja  $NUI = 5/70$  ou 0.07 (7%). Mensagens contendo a palavra **imediatamente** tem uma chance de 0.75 (75%) de serem urgentes, com 0.25 (25%) de chance de serem não urgentes. Em outras palavras, tomando Imediatamente como **I**,  $IU = 15/20$  e  $INU = 5/20$ .

A chance da palavra problema aparecer é  $20/100$ , ou seja 0.2 (20%). No total das 30 mensagens urgentes, 10 continham problema, ou seja  $UP = 10/30$  ou 0.33 (33%). Em contrapartida das 70 mensagens não urgentes, apenas 10/70 continham essa palavra, ou seja  $NUP = 10/70$  ou 0.14 (14%). Mensagens contendo a palavra **Problema** tem uma chance de 0.5 (50%) de serem urgentes, com 0.5 (50%) de chance de não serem urgentes também. Em outras palavras, tomando Problema como **P**,  $PU = 10/20$  e  $PNU = 10/20$ .

A chance da palavra atraso aparecer é  $20/100$ , ou seja 0.2(20%). No total das 30 mensagens urgentes, 8 continham atraso, ou seja  $UA = 8/30$  ou 0.26 (26%). em contrapartida das 70 mensagens não urgentes, 12 mensagens de 70 tinham atraso, ou seja  $NUA = 12/70$  ou 0.17 (17%). Mensagens contendo a palavra **Atraso** tem uma chance de apenas 0.4 (40%) de serem urgentes, com 0.6 (60%) de chance de não serem urgentes. Em outras palavras, tomando Atraso como **A**,  $AU = 8/20$  e  $ANU = 12/20$ .

Essas constatações foram possíveis com uma leitura simples dos dados apresentados.

=====

Correção proposta pela IA:

=====

=====Considerações:

A resposta do estudante demonstra um entendimento básico sobre as probabilidades e como calcular as chances de uma mensagem ser "Urgente" ou "Não Urgente" com base nas palavras presentes. No entanto, há algumas imprecisões e falta de formalismo estatístico na apresentação dos dados.

=====Feedback:

Pontos fortes incluem a identificação correta das probabilidades iniciais para mensagens "Urgentes" e "Não Urgentes", bem como o cálculo de probabilidades condicionais para as palavras "imediatamente", "problema" e "atraso" em relação a essas categorias. No entanto, falta uma clareza na apresentação dos cálculos, especialmente em como as probabilidades condicionais são derivadas. Além disso, as notações estatísticas poderiam ser mais precisas para evitar ambiguidades. Para melhorar, o estudante deve adotar uma abordagem sistemática para apresentar os cálculos, explicitar as fórmulas utilizadas para as probabilidades condicionais e tomar cuidado com a precisão dos cálculos, especialmente ao redor das porcentagens.

=====Pontuação:

2.26089552238806

=====

Correção da Questão 2:

Pergunta: 1b) Uma empresa está desenvolvendo um sistema para classificar mensagens recebidas como "Urgente" ou "Não Urgente" com base nas palavras presentes na mensagem. Foi analisado um conjunto de 100 mensagens, e os dados a seguir foram coletados: Mensagens Urgentes: 30; Mensagens Não Urgentes: 70; Palavra Presente "imediatamente": 15 (Mensagens Urgentes) e 5 (Mensagens Não Urgentes); Palavra Presente "problema": 10 (Mensagens Urgentes) e 10 (Mensagens Não Urgentes); Palavra Presente "atraso": 8 (Mensagens Urgentes) e 12 (Mensagens Não Urgentes). Suponha que uma nova mensagem contenha as palavras "imediatamente" e "problema". Calcule a probabilidade de ser uma mensagem "Urgente" e de ser "Não Urgente" utilizando o teorema de Bayes e classifique a mensagem como "Urgente" ou "Não Urgente".

Diretriz(es):  $\{(P("Urgente"|"imediatamente"&"problema")=(P("urgente")*P("imediatamente"|"urgente")*P("problema"|"urgente"))/(P("problema")*P("imediatamente"))), \text{podendo simplificar para } P("Urgente"|"imediatamente"&"problema")=P("urgente")*P("imediatamente"|"urgente")*P("problema"|"urgente") \sim 0.30 * 0.5 * 0.33 \sim 0.05', 1.0), (P("Classificar o e-mail como \"URGENTE\" baseado no cálculo do máximo argumento entre } P("Urgente"|"imediatamente"&"problema") \text{ e } P("Não Urgente"|"imediatamente"&"problema"), 0.5), (P("Não Urgente"|"imediatamente"&"problema")=(P("Não urgente")*P("imediatamente"|"Não urgente")*P("problema"|"Não urgente"))/(P("problema")*P("imediatamente"))), \text{podendo simplificar para } P("Não Urgente"|"imediatamente"&"problema")=P("Não urgente")*P("imediatamente"|"Não urgente")*P("problema"|"Não urgente") \sim 0.70 * 0.07 * 0.014 \sim 0.001', 1.0)\}$

=====Resposta do estudante da questão 2:

A2)  $P(U|■■■,■P■) = P(■■■|U) P(■P■|U) P(U) / P(■■■) P(■P■)$

$P(U|■■■,■P■) = 0.5 \times 0.33 \times 0.3 / 0.2 \times 0.2$

$P(U|■■■, ■P■) = 1.2375$

Urgente = 1.2375

$P(NU|■■■,■P■) = P(■■■|NU) P(■P■|NU) P(NU) / P(■■■) P(■P■)$

$P(NU|■■■,■P■) = 0.25 \times 0.5 \times 0.7 / 0.2 \times 0.2$

$P(NU|■■■, ■P■) = 2.1875$

Não urgente = 2.1875

=====

Correção proposta pela IA:

=====

=====Considerações:

A resposta do estudante apresenta alguns erros nos cálculos e na aplicação do teorema de Bayes. Primeiramente, é necessário calcular as probabilidades condicionais corretamente. A probabilidade de uma mensagem ser urgente dado que contém as palavras "imediatamente" e "problema" pode ser calculada como  $P(U|I,P) = P(U) * P(I|U) * P(P|U) / P(I) * P(P)$ .

=====Feedback:

Os principais pontos a serem corrigidos incluem o cálculo das probabilidades condicionais e a aplicação do teorema de Bayes. Além disso, os valores calculados precisam ser interpretados corretamente para classificar a mensagem como "Urgente" ou "Não Urgente". Para atingir a nota máxima, o estudante deve demonstrar uma compreensão clara do teorema de Bayes e como aplicá-lo corretamente a esse problema.

=====Pontuação:

0.6376146788990825

=====

Correção da Questão 3:

Pergunta: 2a) Árvores de decisão são modelos amplamente utilizados para classificação e regressão. Explique como o conceito de ganho de informação é utilizado na construção de uma árvore de decisão. Detalhe mostrando o uso do cálculo de entropia e ganho de informação em um problema hipotético.

Diretriz(es):  $\{(P(\text{"Demonstrar o cálculo de entropia final como sendo } H(\text{"depois da partição"})=P(\text{"amostras irem para conjunto 1"})*H(\text{"conjunto 1 após partição"})+P(\text{"amostras irem para conjunto 2"})*H(\text{"conjunto 2 após partição"}), 1.0), (P(\text{"Explicar que o conceito de Information Gain (IG) como sendo } IG=H(\text{"antes da partição"})-H(\text{"depois da partição"}) \text{ é essencial$

para decidir qual atributo/partição escolher a cada nível da árvore', 1.0), ('Demonstrar o cálculo de entropia inicial como sendo  $H(\text{"antes da partição"}) = -p(\text{"classe1"}) \cdot \log_2(P(\text{"classe1"})) - \dots - p(\text{"classeN"}) \cdot \log_2(P(\text{"classeN"}))$ ', 1.0))  
=====Resposta do estudante da questão 3:

B1) O ganho de informação serve para comparar dados na premissa de definir se a sua inclusão é válida ou não por virtude da similaridade com outros dados. Por exemplo: Se em uma tabela qualquer temos as colunas ■Ventos fortes■ e ■Chuva■ sendo associadas à linha de ■Evento ocorreu■ com uma classificação binária de ■Sim■ ou ■Não■ é válido verificar se há algum ganho por usar de ambas ou se elas são indicadores redundantes por sua similaridade. Por mérito de exemplo, consideremos que na coluna de ■Ventos fortes■, temos 6 sim e 4 não quando relacionado à ■Evento ocorreu■, e da mesma forma em ■Chuva■ temos 5 sim e 5 não. Fazendo o cálculo da entropia para encontrar o ganho:

$$(-0.6 \log_2(0.6)) = 0.4421$$

$$(-0.5 \log_2(0.5)) = 0.5$$

$$0.4421 - 0.5 = 0.0579$$

Pela proximidade do resultado à 0, podemos constatar que ■Ventos fortes■ e ■Chuva■ tem um ganho ruim, e seria viável ao sistema remover uma dessas variáveis.

=====

Correção proposta pela IA:

=====

=====Considerações:

A resposta do estudante demonstra um entendimento básico do conceito de ganho de informação e sua aplicação em árvores de decisão. No entanto, falta uma explicação clara e detalhada sobre como o cálculo de entropia é utilizado para determinar o ganho de informação. Além disso, o exemplo fornecido não ilustra corretamente o cálculo da entropia inicial e final, nem a fórmula para calcular o ganho de informação.

=====Feedback:

Para melhorar a resposta, o estudante deve começar explicando o cálculo da entropia inicial, utilizando a fórmula  $H(\text{"antes da partição"}) = -p(\text{"classe1"}) \cdot \log_2(P(\text{"classe1"})) - \dots - p(\text{"classeN"}) \cdot \log_2(P(\text{"classeN"}))$ . Em seguida, deve ilustrar o cálculo da entropia final, após a partição, utilizando a fórmula  $H(\text{"depois da partição"}) = P(\text{"amostras irem para conjunto 1"}) \cdot H(\text{"conjunto 1 após partição"}) + P(\text{"amostras irem para conjunto 2"}) \cdot H(\text{"conjunto 2 após partição"})$ . Finalmente, o estudante deve aplicar a fórmula do ganho de informação,  $IG = H(\text{"antes da partição"}) - H(\text{"depois da partição"})$ , para demonstrar como escolher o atributo ou partição mais adequados a cada nível da árvore de decisão.

=====Pontuação:

1.125

=====

Correção da Questão 4:

Pergunta: 2b) Árvores de decisão são modelos amplamente utilizados para classificação e regressão. Suponha que você está desenvolvendo um sistema de suporte à decisão, aponte dois critérios de parada que você utilizaria na construção da árvore de decisão para garantir um modelo eficiente?

Diretriz(es): {'Apontar a limitação da quantidade de níveis da árvore (profundidade máxima atingida) como um possível critério de parada.', 1.25), ('Apontar qualquer outra limitação plausível para que uma árvore de decisão não continue a realizar partições com base em um número de amostras presentes em um dado nó.', 1.25), ('Apontar a limitação de número insuficiente de amostras no nó para realizar nova partição como um possível critério de parada.', 1.25), ('Apontar a limitação de entropia das classes (labels) atingir zero no nó como um possível critério de parada.', 1.25)}

=====Resposta do estudante da questão 4:

B2) De forma sucinta, o primeiro seria por comprimento da árvore, limitando o seu número de passos por um valor arbitrário para impedir uma execução infinita. O outro seria por ■pureza■, ou seja, o quão próximo de 0 se encontra a entropia, deixando que o processo seja executado até que se alcance uma entropia mínima desejada sem deixar que continue com uma execução desnecessária.

=====

Correção proposta pela IA:

=====

=====Considerações:

A resposta do estudante aborda dois critérios de parada importantes na construção de árvores de decisão, que são a limitação do comprimento da árvore (profundidade máxima atingida) e a pureza do nó (entropia mínima desejada). No entanto, a resposta poderia ser mais precisa em relação às diretrizes fornecidas, especialmente quanto à nomenclatura e ao detalhamento dos critérios.

=====Feedback:

Um ponto forte da resposta é a menção à limitação do comprimento da árvore, o que se alinha com a diretriz de limitar a quantidade de níveis da árvore. Além disso, a menção à pureza do nó como critério de parada toca na ideia de limitar a entropia das classes, o que é similar à diretriz de apontar a limitação de entropia atingir zero no nó. No entanto, para melhorar, o estudante poderia explicitar mais claramente esses conceitos, usando terminologias como "profundidade máxima" e "entropia mínima" ou "número insuficiente de amostras" para justificar a parada. Além disso, seria útil mencionar que a limitação do número de amostras no nó é outro critério válido para parar a árvore, o que está de acordo com as diretrizes.

=====Pontuação:

2.0

=====