

Nome: LUÃ MOREIRA PONCIANO

Pontuação Total da Avaliação: 7.80 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".

=====Rubrica(s)

('Acertar a resposta aproximada de $P(\text{"não urgente"}) = 70/100 = 0.70$ ', 0.32)

('Acertar a resposta aproximada de $P(\text{"problema"}|\text{"não urgente"}) = 10/70 \sim 0.14$ ', 0.32)

('Acertar a resposta aproximada de $P(\text{"imediatamente"}|\text{"não urgente"}) = 5/70 \sim 0.07$ ', 0.32)

('Acertar a resposta aproximada de $P(\text{"problema"}|\text{"urgente"}) = 10/30 \sim 0.33$ ', 0.32)

('Acertar a resposta aproximada de $P(\text{"atraso"}|\text{"urgente"}) = 8/30 \sim 0.27$ ', 0.32)

('Acertar a resposta aproximada de $P(\text{"urgente"}) = 30/100 = 0.30$ ', 0.32)

('Acertar a resposta aproximada de $P(\text{"imediatamente"}|\text{"urgente"}) = 15/30 = 0.5$ ', 0.32)

('Acertar a resposta aproximada de $P(\text{"atraso"}|\text{"não urgente"}) = 12/70 \sim 0.17$ ', 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, $U = 70/100$ e $NU = 30/100$.

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 apresenta boa compreensão sobre as probabilidades das mensagens serem "Urgentes" ou "Não Urgentes" com base nos dados fornecidos. Ele corretamente calculou as probabilidades de uma mensagem ser "Urgente" ($30/100 = 0.30$) e "Não Urgente" ($70/100 = 0.70$). Além disso, o estudante calculou as probabilidades condicionais para cada palavra em relação às mensagens "Urgentes" e "Não Urgentes", demonstrando entendimento das relações entre as palavras e a classificação das mensagens.

=====Feedback:

Pontos fortes incluem o cálculo correto das probabilidades básicas e a compreensão das probabilidades condicionais. No entanto, para melhorias, é sugerido que o estudante detalhe

melhor os passos de cálculo e considere discutir brevemente a implicação desses cálculos para o sistema de classificação de mensagens. Além disso, uma discussão mais aprofundada sobre como essas probabilidades condicionais poderiam ajudar a melhorar o sistema de classificação seria benéfica.

=====Pontuação:

2.4

=====

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".

=====Rubrica(s)

('Calcular a probabilidade $P(\text{"Não Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = (P(\text{"Não urgente"}) * P(\text{"imediatamente"}|\text{"Não urgente"}) * P(\text{"problema"}|\text{"Não urgente"})) / (P(\text{"problema"}) * P(\text{"imediatamente"}))$, podendo simplificar para $P(\text{"Não Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = P(\text{"Não urgente"}) * P(\text{"imediatamente"}|\text{"Não urgente"}) * P(\text{"problema"}|\text{"Não urgente"}) \sim 0.70 * 0.07 * 0.014 \sim 0.001$, 1.0)

('Calcular a probabilidade $P(\text{"Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = (P(\text{"urgente"}) * P(\text{"imediatamente"}|\text{"urgente"}) * P(\text{"problema"}|\text{"urgente"})) / (P(\text{"problema"}) * P(\text{"imediatamente"}))$, podendo simplificar para $P(\text{"Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = P(\text{"urgente"}) * P(\text{"imediatamente"}|\text{"urgente"}) * P(\text{"problema"}|\text{"urgente"}) \sim 0.30 * 0.5 * 0.33 \sim 0.05$, 1.0)

('Classificar o e-mail como "URGENTE" baseado no cálculo do máximo argumento entre $P(\text{"Urgente"}|\text{"imediatamente"}\&\text{"problema"})$ e $P(\text{"Não Urgente"}|\text{"imediatamente"}\&\text{"problema"})$, 0.5)

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

A2) $P(U|I, "P") = P(I|U) P(P|U) P(U) / P(I) P("U")$

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

$P(U|I, "P") = 1.2375$

Urgente = 1.2375

$P(NU|I, "P") = P(I|NU) P(P|NU) P(NU) / P(I) P("P")$

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

$P(NU|I, "P") = 2.1875$

NÃO urgente = 2.1875

=====

Correção proposta pela IA:

=====

=====Considerações:

A resposta do estudante contém alguns erros de cálculo e não segue corretamente as rubricas fornecidas. No entanto, o estudante demonstrou um conhecimento básico sobre o assunto, tentando aplicar o teorema de Bayes para calcular as probabilidades. Para calcular a probabilidade de uma mensagem ser "Urgente" dado que contenha as palavras "imediatamente" e "problema", o estudante deveria ter calculado $P(\text{"Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = (P(\text{"Urgente"}) * P(\text{"imediatamente"}|\text{"Urgente"}) * P(\text{"problema"}|\text{"Urgente"})) / (P(\text{"problema"}) * P(\text{"imediatamente"}))$.

=====Feedback:

O estudante cometeu erros nos cálculos. Primeiro, para calcular as probabilidades condicionais, é necessário entender que $P(\text{"imediatamente"}|\text{"Urgente"}) = 15/30 = 0.5$, $P(\text{"problema"}|\text{"Urgente"}) = 10/30 = 0.33$, $P(\text{"imediatamente"}|\text{"Não Urgente"}) = 5/70 = 0.07$ e $P(\text{"problema"}|\text{"Não Urgente"}) = 10/70 = 0.14$. Além disso, para calcular $P(\text{"problema"})$ e $P(\text{"imediatamente"})$, precisamos considerar todas as ocorrências, então $P(\text{"problema"}) = (10+10)/100 = 0.2$ e $P(\text{"imediatamente"}) = (15+5)/100 = 0.2$. Portanto, $P(\text{"Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = (0.3 * 0.5 * 0.33) / (0.2 * 0.2) = 0.2475$. Para $P(\text{"Não Urgente"}|\text{"imediatamente"}\&\text{"problema"})$, temos $P(\text{"Não Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = (0.7 * 0.07 * 0.14) / (0.2 * 0.2) = 0.016225$. Como $0.2475 >$

0.016225, a mensagem deve ser classificada como "Urgente".

=====Pontuação:

1.7

=====

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.

=====Rubrica(s)

('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)

('Explicar que o conceito de Information Gain (IG) como sendo $IG = H(\text{"antes da partição"}) - H(\text{"depois da partição"})$ é essencial para decidir qual atributo/partição escolher a cada nível da árvore', 1.0)

('Demonstrar o cálculo de entropia inicial como sendo $H(\text{"antes da partição"}) = -p(\text{"classe1"}) * \log_2(P(\text{"classe1"})) - \dots - p(\text{"classeN"}) * \log_2(P(\text{"classeN"}))$ ', 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 motivo 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 aborda o conceito de ganho de informação e sua utilidade na construção de árvores de decisão, embora de forma simplificada e com alguns erros conceituais. Ele tenta calcular a entropia para encontrar o ganho de informação, mas não aplica corretamente a fórmula de entropia inicial e final. Além disso, falta uma explicação clara sobre como o conceito de Information Gain (IG) é essencial para decidir qual atributo ou partição escolher a cada nível da árvore. A aplicação do cálculo de entropia e ganho de informação em um problema hipotético é iniciada, mas não levada a cabo de forma consistente com as rubricas.

=====Feedback:

Para melhorar, é fundamental que o estudante compreenda e aplique corretamente a fórmula de entropia inicial $H(\text{"antes da partição"}) = -p(\text{"classe1"}) * \log_2(P(\text{"classe1"})) - \dots - p(\text{"classeN"}) * \log_2(P(\text{"classeN"}))$, e então calcule a entropia final $H(\text{"depois da partição"}) = P(\text{"amostras irem para conjunto 1"}) * H(\text{"conjunto 1 após partição"}) + P(\text{"amostras irem para conjunto 2"}) * H(\text{"conjunto 2 após partição"})$. Além disso, deve explicar claramente como o conceito de $IG = H(\text{"antes da partição"}) - H(\text{"depois da partição"})$ é crucial para a escolha do atributo ou partição mais informativa a cada nodo da árvore. Um exemplo detalhado, aplicando esses conceitos a um problema hipotético, ajudaria a solidificar a compreensão do estudante sobre o tema.

=====Pontuação:

1.2

=====

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?

=====Rubrica(s)

('Mencionar e/ou explicar dois apenas um critério de parada válido na construção de árvores de decisão.', 1.25)

('Mencionar e/ou explicar dois ou mais critérios de parada válidos na construção de árvores de decisão.', 2.5)

('Não mencionar e/ou explicar qualquer critério de parada na construção de árvores de decisão.', 0.0)

=====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 quanto 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 menciona dois critérios de parada para a construção de árvores de decisão, que são o comprimento da árvore (limitando o número de passos) e a pureza (ou seja, a entropia). Embora a resposta seja concisa e aponte na direção certa, ela poderia ser mais clara e detalhada, especialmente para um estudante com conhecimento básico sobre o assunto.

=====Feedback:

Pontos fortes da resposta incluem a menção explícita de dois critérios de parada, o que já se alinha parcialmente com o esperado pela rubrica. No entanto, para melhorar e atingir a nota máxima, seria importante explicar mais detalhadamente cada critério, por exemplo, como o comprimento da árvore afeta a complexidade do modelo e como a pureza (entropia) impacta a precisão da classificação ou regressão. Além disso, exemplos práticos ou uma breve discussão sobre como esses critérios podem ser implementados em algoritmos de árvore de decisão seriam valiosos.

=====Pontuação:

2.5

=====