

Nome: LUÃ MOREIRA PONCIANO

Pontuação Total da Avaliação: 7.90 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 $P(\text{"imediatamente"}|\text{"urgente"})=15/30=0.5'$, 0.32)

('Acertar a resposta aproximada de $P(\text{"problema"}|\text{"urgente"})=10/30\sim0.33'$, 0.32)

('Acertar a resposta aproximada de $P(\text{"atraso"}|\text{"urgente"})=8/30\sim0.27'$, 0.32)

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

('Acertar a resposta aproximada de $P(\text{"problema"}|\text{"não urgente"})=10/70\sim0.14'$, 0.32)

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

('Acertar a resposta aproximada de $P(\text{"atraso"}|\text{"não urgente"})=12/70\sim0.17'$, 0.32)

('Acertar a resposta aproximada de $P(\text{"imediatamente"}|\text{"não urgente"})=5/70\sim0.07'$, 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 demonstra uma boa compreensão dos conceitos básicos de probabilidade e condicionalidade. Ele calculou corretamente as probabilidades de uma mensagem ser "Urgente" ($U=30/100=0.3$) e "Não Urgente" ($NU=70/100=0.7$). Além disso, o estudante corretamente identificou as probabilidades condicionais para cada palavra ("imediatamente", "problema" e "atraso") em relação às mensagens "Urgentes" e "Não Urgentes". No entanto, é importante notar que a apresentação poderia ser mais clara e organizada, e algumas pequenas imprecisões foram encontradas na descrição das probabilidades.

=====Feedback:

Para alcançar a nota máxima, o estudante poderia ter detalhado melhor cada passo do cálculo,

especialmente para as probabilidades condicionais. Além disso, uma apresentação mais sistemática dos resultados, possivelmente em formato de tabela ou com uma clara separação entre os cálculos, ajudaria na compreensão. Outro ponto é a necessidade de explicitar melhor as fórmulas ou métodos utilizados para calcular as probabilidades, especialmente para os leitores não familiarizados com esses conceitos. No geral, a resposta está próxima do esperado, mas com ajustes na organização e detalhamento, poderia atingir a perfeição.

=====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)

('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)
('Calcular a probabilidade $P(\text{"Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = (P(\text{"urgente"})*P(\text{"imediatamente"}|\text{"urgente"})*P(\text{"problema"}|\text{"urgente"})) / (P(\text{"problema"})*P(\text{"imediatamente"}|\text{"problema"}))$, 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)

('Calcular a probabilidade $P(\text{"Não Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = (P(\text{"Não urgente"}) * P(\text{"imediatamente"}|\text{"Não urgente"}) * P(\text{"problema"}|\text{"Não urgente"})) / (P(\text{"problema"})*P(\text{"imediatamente"}|\text{"problema"}))$, 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)

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

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

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

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

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 apresenta uma tentativa de aplicar o teorema de Bayes para calcular as probabilidades de uma mensagem ser urgente ou não urgente com base nas palavras presentes. No entanto, há erros significativos nos cálculos e na aplicação das fórmulas.

=====Feedback:

O estudante tentou calcular $P(U|I, "P")$ e $P(NU|I, "P")$, mas cometeu erros nos passos intermediários. Para calcular corretamente, precisamos aplicar o teorema de Bayes da forma correta. Primeiro, calcule as probabilidades condicionais e a priori: $P(I|U) = 15/30 = 0.5$, $P(P|U) = 10/30 = 0.33$, $P(U) = 30/100 = 0.3$, $P(I)$ precisa ser calculada considerando todos os casos em que "I" aparece, então $P(I) = (15+5)/100 = 0.2$. De forma similar, $P(P) = (10+10)/100 = 0.2$. Para $P(NU|I, "P")$, precisamos de $P(I|NU) = 5/70$, $P(P|NU) = 10/70$ e $P(NU) = 70/100$. Aplicando o teorema de Bayes corretamente, $P(U|I, "P") = (P(I|U) * P(P|U) * P(U)) / (P(I) * P("P"))$ e $P(NU|I, "P") = (P(I|NU) * P(P|NU) * P(NU)) / (P(I) * P("P"))$. Com esses valores, podemos calcular as probabilidades corretas e compará-las para decidir se a mensagem é urgente ou não.

=====Pontuação:

1.5

=====

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)

('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)

('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)

=====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 demonstra uma compreensão básica do conceito de ganho de informação e sua aplicação na construção de árvores de decisão. Ele fornece um exemplo hipotético e tenta aplicar o cálculo de entropia para encontrar o ganho de informação. No entanto, a resposta carece de clareza e precisão em sua explicação, especialmente ao calcular a entropia final como sendo $H(\text{"depois da partição"})$ e ao demonstrar a importância do conceito de Information Gain (IG) como sendo $IG = H(\text{"antes da partição"}) - H(\text{"depois da partição"})$.

=====Feedback:

Para melhorar, o estudante deve trabalhar em explicar de forma clara e detalhada como o cálculo de entropia é aplicado antes e depois da partição, usando as fórmulas corretas, $H(\text{"antes da partição"}) = -p(\text{"classe1"}) * \log_2(P(\text{"classe1"})) - \dots - p(\text{"classeN"}) * \log_2(P(\text{"classeN"}))$ e $H(\text{"depois da partição"}) = P(\text{"amostras irem para conjunto 1"}) * H(\text{"conjunto 1 após partição"}) + P(\text{"amostras irem para conjunto 2"}) * H(\text{"conjunto 2 após partição"})$. Além disso, deve-se enfatizar a importância do Information Gain (IG) na escolha do atributo mais relevante para a partição em cada nível da árvore de decisão.

=====Pontuação:

1.5

=====

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 numero 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 apresenta uma boa compreensão dos critérios de parada na construção de árvores de decisão. Ele menciona dois critérios: o comprimento da árvore e a pureza (ou entropia). No entanto, a explicação poderia ser mais detalhada e clara. O estudante também poderia ter fornecido exemplos ou justificativas para apoiar sua resposta.

=====Feedback:

Para melhorar a resposta, o estudante poderia fornecer mais detalhes sobre como esses critérios são utilizados na prática. Por exemplo, ele poderia explicar como o comprimento da árvore afeta a complexidade do modelo e como a pureza é calculada e utilizada como critério de parada. Além disso, o estudante poderia mencionar outros critérios de parada, como o número mínimo de amostras por nó ou o número máximo de divisões. Isso demonstraria uma compreensão mais profunda do assunto e permitiria que o estudante fornecesse uma resposta mais completa e precisa.

=====Pontuação:

2.5

=====