

Nome: Francisco Lucas Benvindo

Pontuação Total da Avaliação: 6.76 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):  $\{ (P("urgente")=30/100=0.30, 0.32), (P("atraso"|"não urgente")=12/70\sim0.17', 0.32), (P("imediatamente"|"urgente")=15/30=0.5', 0.32), (P("imediatamente"|"não urgente")=5/70\sim0.07', 0.32), (P("atraso"|"urgente")=8/30\sim0.27', 0.32), (P("problema"|"urgente")=10/30\sim0.33', 0.32), (P("não urgente")=70/100=0.70', 0.32), (P("problema"|"não urgente")=10/70\sim0.14', 0.32) \}$

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

1A)  $P(\hat{\text{Urgente}}) = 30/100 = 0.3$   
 $P(\hat{\text{Não urgente}}) = 70/100 = 0.7$   
 $P(\hat{\text{imediatamente}} | \hat{\text{Urgente}}) = 15/30 = 0.5$   
 $P(\hat{\text{problema}} | \hat{\text{Urgente}}) = 10/30 = 0.33$   
 $P(\hat{\text{atraso}} | \hat{\text{Urgente}}) = 8/30 = 0.27$   
 $P(\hat{\text{imediatamente}} | \hat{\text{Não urgente}}) = 5/70 = 0.07$   
 $P(\hat{\text{problema}} | \hat{\text{Não urgente}}) = 10/70 = 0.14$   
 $P(\hat{\text{atraso}} | \hat{\text{Não urgente}}) = 12/70 = 0.17$

=====

Correção proposta pela IA:

=====

=====Considerações:

A resposta do estudante apresenta alguns erros e inconsistências em relação às diretrizes fornecidas. Primeiramente, as probabilidades de uma mensagem ser "Urgente" e "Não Urgente" foram calculadas corretamente como  $P("Urgente") = 0,3$  e  $P("Não Urgente") = 0,7$ , o que está de acordo com as diretrizes.

=====Feedback:

No entanto, os cálculos das probabilidades condicionais apresentam erros. Por exemplo,  $P("imediatamente" | "Urgente")$  foi calculado como  $15/20 = 0,75$ , mas deveria ser  $15/30 = 0,5$ , considerando que há 30 mensagens urgentes no total. Da mesma forma,  $P("imediatamente" | "Não Urgente")$  foi calculado como  $5/20$ , mas deveria ser  $5/70 = 0,0714$ , considerando que há 70 mensagens não urgentes. As probabilidades condicionais para "problema" e "atraso" também precisam ser recalculadas com base nos totais corretos de mensagens urgentes e não urgentes.

=====Pontuação:

1.92

=====

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):  $\{ (P("Urgente"|"imediatamente" \& "problema") = (P("urgente") * P("imediatamente"|"urgente") * P("problema"|"urgente")) / (P("problema") * P("imediatamente"|"urgente") + P("não urgente") * P("imediatamente"|"urgente") * P("problema"|"urgente")) \}$ , podendo simplificar para  $P("Urgente"|"imediatamente" \& "problema") = P("urgente") * P("imediatamente"|"urgente") * P("problema"|"urgente") / (P("problema") * (P("imediatamente"|"urgente") * (P("urgente") + P("não urgente") * P("imediatamente"|"urgente") * P("problema"|"urgente"))))$

nte")  $\sim 0.30 * 0.5 * 0.33 \sim 0.05$ , 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), ( $P(\text{"Não Urgente"}|\text{"imediatamente"}\&\text{"problema"}) = (P(\text{"Não urgente"}) * P(\text{"imediatamente"}|\text{"Não urgente"}) * P(\text{"problema"}|\text{"Não urgente"})) / (P(\text{"problema"}) * P(\text{"imediatamente"}))$ ), 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)}

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

1B)  $P(\text{"Urgente"} | \text{"imediatamente"} \& \text{"problema"}) = P(\text{"Urgente"}) * P(\text{"imediatamente"} | \text{"Urgente"}) * P(\text{"problema"} | \text{"Urgente"})$

Substituindo:  $0.3 * 0.75 * 0.5 = 0.1125$

$P(\text{"Não Urgente"} | \text{"imediatamente"} \& \text{"problema"}) =$

$P(\text{"Não Urgente"}) * P(\text{"imediatamente"} | \text{"Não Urgente"}) * P(\text{"problema"} | \text{"Não Urgente"})$

Substituindo:  $0.7 * 0.25 * 0.5 = 0.0875$

A mensagem pode ser classificada como urgente.

=====Correção proposta pela IA:

=====Considerações:

A resposta do estudante apresenta uma abordagem inicial correta ao aplicar o teorema de Bayes para calcular as probabilidades de uma mensagem ser "Urgente" ou "Não Urgente" com base nas palavras presentes. No entanto, é necessário considerar a probabilidade a priori e a probabilidade condicional de maneira mais precisa, incluindo a consideração do denominador do teorema de Bayes, que normaliza as probabilidades.

=====Feedback:

O estudante começou aplicando o teorema de Bayes de forma simplificada para calcular  $P(\text{"Urgente"} | \text{"imediatamente"} \& \text{"problema"})$  e  $P(\text{"Não Urgente"} | \text{"imediatamente"} \& \text{"problema"})$ . No entanto, para uma resposta completa, seria necessário calcular corretamente as probabilidades a priori e condicionais com base nos dados fornecidos. As probabilidades a priori são  $P(\text{"Urgente"}) = 30/100 = 0,3$  e  $P(\text{"Não Urgente"}) = 70/100 = 0,7$ . As probabilidades condicionais necessárias são  $P(\text{"imediatamente"} | \text{"Urgente"})$ ,  $P(\text{"problema"} | \text{"Urgente"})$ ,  $P(\text{"imediatamente"} | \text{"Não Urgente"})$  e  $P(\text{"problema"} | \text{"Não Urgente"})$ , que podem ser calculadas com base nos dados fornecidos. Além disso, para uma comparação completa, seria ideal calcular o denominador comum para ambas as probabilidades condicionais, considerando as probabilidades de encontrar as palavras independentemente da urgência. Para melhorias, é essencial revisar a aplicação do teorema de Bayes, considerando todas as partes necessárias para uma análise precisa.

=====Pontuação:

1.2

=====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): {'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:

2A) O conceito de ganho de informação é utilizado em uma árvore de decisão para reduzir a entropia, ou grau de confusão nos dados. Ao selecionar determinada variável e realizar uma partição com nela, o ideal é que haja um ganho de informação e redução da entropia. Exemplo:

Dado um conjunto de dados tabular que registra a ocorrência de idas a praia com 12 linhas, e com as variáveis independentes de temperatura, clima, força do vento, estado da maré entre outras e uma variável alvo "ida" que pode ser classificada como "foi" ou "não foi". Considere que dos 12 registros, 7 são classificados como "foi" e 5 como "não foi". O cálculo de entropia antes da tomada de qualquer decisão se da por:

Legenda:

f = foi

nf = não foi

$$-P(f) \cdot \log_2(P(f)) - P(nf) \cdot \log_2(P(nf)) = -7/12 \cdot \log_2(7/12) - 5/12 \cdot \log_2(5/12) = 0.455 + 0.527 = 0.982$$

Dado que por exemplo seleciona-se a variável de temperatura e cria-se uma partição com a condição temperatura > 28. Realiza-se novamente o cálculo de entropia, onde a temperatura foi maior que 28 e a ida foi classificada como "foi" e também como "não foi" o mesmo ocorre para quando a temperatura foi menor que 28. Ao final, somam-se as duas entropias e subtrai-se da entropia antes da tomada de qualquer decisão, calculando assim o ganho de informação

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. Ele explica corretamente que o objetivo é reduzir a entropia nos dados ao selecionar uma variável e realizar uma partição. Além disso, o estudante fornece um exemplo hipotético para ilustrar o cálculo de entropia antes e depois da partição.

Feedback:

No entanto, a resposta do estudante carece de detalhes importantes para atingir a nota máxima. Em primeiro lugar, ele não calcula explicitamente a entropia após a partição e não apresenta a fórmula correta para o cálculo de entropia final como  $H(\text{"depois da partição"}) = P(\text{"amostras irem para conjunto 1"}) \cdot H(\text{"conjunto 1 após partição"}) + P(\text{"amostras irem para conjunto 2"}) \cdot H(\text{"conjunto 2 após partição"})$ . Além disso, não enfatiza claramente que o conceito de Information Gain (IG) é essencial para decidir qual atributo ou partição escolher a cada nível da árvore, utilizando a fórmula  $IG = H(\text{"antes da partição"}) - H(\text{"depois da partição"})$ .

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?

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:

2B) Utilizaria os critérios de parada de nível da árvore, determinando um nível máximo onde a árvore deve parar de expandir, evitando com que a árvore de decisão fique muito extensa e o critério da pureza ou homogeneidade dos dados, onde os dados de determinado nível são da mesma classe.

Correção proposta pela IA:

Considerações:

A resposta do estudante aborda dois critérios de parada importantes na construção de árvores de decisão, que são a profundidade máxima da árvore e a pureza ou homogeneidade dos dados em um

nó. No entanto, a resposta pode ser aprimorada para atender às diretrizes fornecidas.

=====Feedback:

Pontos fortes da resposta incluem a menção à limitação da profundidade da árvore (nível máximo), o que está alinhado com as diretrizes. Entretanto, a resposta poderia ser melhorada explicitando critérios de parada mais detalhados, como o número insuficiente de amostras no nó ou a entropia das classes atingir zero, para garantir que a árvore de decisão seja eficiente e evite sobreajuste.

=====Pontuação:

2.1367521367521367

=====