

Etapas da CRISP-DM:	
1. Compreensão do Negócio:	Identificar objetivos do projeto e converter em metas específicas de mineração de dados.
2. Compreensão dos Dados:	Coleta, descrição e análise inicial dos dados.
3. Preparação dos Dados:	Limpeza, integração e transformação dos dados para o formato adequado.
4. Modelagem:	Seleção e aplicação de técnicas de modelagem de dados.
5. Avaliação:	Verificar se o modelo atende aos objetivos estabelecidos.
6. Implementação:	Aplicar os modelos em ambiente operacional para uso contínuo.

Nº _____ CURSO _____

NOME eduardovski**GRUPO 1**

(4 valores)

QUESTÃO 1

RESPONDA ÀS Q

No desenvolvimento de sistemas de aprendizagem automática (*machine learning*) podem ser utilizados diferentes paradigmas de aprendizagem.

Neste contexto pretende-se que:

- caracterize os paradigmas de aprendizagem supervisionada, não supervisionada e por reforço;
- apresente dois exemplos de técnicas de cada paradigma, ilustrando-os com casos de aplicação.

QUESTÃO 2

O processo de desenvolvimento de uma solução de aprendizagem automática envolve diversas etapas, que podem diferir de acordo com a metodologia escolhida.

Tendo em consideração a metodologia CRISP-DM, pretende-se que enumere e descreva as suas etapas.

GRUPO 2

(4 valores)

Responda às questões deste grupo no espaço reservado PREENCHENDO OS ESPAÇOS VAZIOS com as expressões devidas de modo que a afirmação seja correta.

QUESTÃO 1

No contexto da utilização de técnicas de aprendizagem automática (*machine learning*), a adoção de uma metodologia para a extração de conhecimento descreve e cria _____ **etapas** _____ pelos quais deverá passar o desenvolvimento de um projeto de extração de conhecimento para _____ **a tomada de decisao** _____.

QUESTÃO 2

A metodologia de extração de conhecimento que se desenvolve em 5 etapas, a saber,

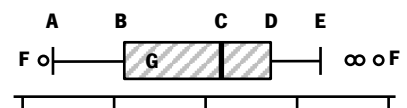
_____ **sample / amostra** _____, _____ **Explore** _____, _____ **Modify** _____, _____ **Model** _____ e _____ **Acess / avaliacao** _____, denomina-se SEMMA.

QUESTÃO 3

Máquina de Vetores de Suporte (*Support Vector Machine*) é uma técnica _____ **supervisionada** _____ de aprendizagem automática que pode ser utilizada para resolver problemas de _____ **regressao** _____ e de _____ **classificacao** _____.

QUESTÃO 4

Num diagrama de caixa (*boxplot*), como no exemplo à direita, o ponto **C** corresponde à _____ **mediana** _____, a caixa **G** representa _____ **intervalo interquartil** _____ dos dados do estudo, e os círculos **F** identificam os valores _____ **outliers** _____ do *dataset*.



GRUPO 3

(6 valores)

RESPONDA ÀS QUESTÕES DESTE GRUPO NO ESPAÇO RESERVADO.

Considere o *dataset Titanic*, utilizado por diversas vezes ao longo do semestre. Considere também o excerto de código apresentado na Figura 1, onde é apresentada a construção e avaliação de um modelo de aprendizagem automática.

QUESTÃO 1

O excerto apresentado contém imprecisões. Identifique e corrija-as utilizando o espaço disponível ao lado da figura (não deve copiar todo o excerto, mas apenas aquilo que corrigiu).

```
[1] df = pd.read_csv('titanic_dataset.csv')
[2] x = df.drop(['Survived', 'Age', 'PassengerId', 'Name',
               'Ticket', 'Cabin', 'Embarked', 'Sex'], axis=1)
[3] y = df['Survived']
[4] sex_ohe = pd.merge(df['Sex'], drop_first=True)
[5] embarked_ohe = pd.merge(df['Embarked'], drop_first=True)
[6] X = pd.concat([X, sex_ohe, embarked_ohe], axis=1)
[7] X_train, X_test, y_train, y_test =
    train_test_split(y, X, test_size=0.3)
[8] model = Sequential()
[9] model.add(Dense(16, input_dim=y.shape[1],
                  activation='relu'))
[10] model.add(Dense(8, activation='relu'))
[11] model.add(Dense(1, activation='sigmoid'))
[12] model.compile(loss = 'binary_crossentropy',
                  optimizer = 'adam',
                  metrics = ['mse'])
[13] model.transform(X_train, y_train, epochs=50,
                   batch_size=32)
[14] loss, acc = model.evaluate(X_train, y_train)
```

python

```
[4] sex_ohe = pd.get_dummies(df['Sex'], drop_first=True) # Substituir pd.me
[5] embarked_ohe = pd.get_dummies(df['Embarked'], drop_first=True) # Correçã
[7] X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
[9] model.add(Dense(16, input_dim=X.shape[1], activation='relu')) # Corrigir
[13] model.fit(X_train, y_train, epochs=50, batch_size=32) # Substituir 'tra
```

Figura 1. Excerto de um modelo de aprendizagem.

QUESTÃO 2

Identifique a técnica de aprendizagem utilizada no excerto de código apresentado na Figura 1, e indique quatro hiperparâmetros passíveis de serem modificados para afinar o modelo.

Questão 2: Técnica de aprendizagem e hiperparâmetros

- **Técnica de aprendizagem:** A técnica usada é **Redes Neurais Artificiais (RNA)**, uma abordagem de aprendizagem supervisionada.
- **Quatro hiperparâmetros ajustáveis:**
 1. **Número de neurónios em cada camada:** Controla a capacidade de aprendizado da rede.
 2. **Taxa de aprendizado:** Define o passo dado durante a otimização.
 3. **Função de ativação:** Pode variar (e.g., ReLU, sigmoid, tanh) para diferentes comportamentos do modelo.
 4. **Número de épocas (epochs):** Define quantas vezes o modelo passa pelos dados durante o treinamento.

QUESTÃO 3

Admita que o *dataset Titanic* não está balanceado. Descreva de que forma este desbalanceamento influencia o modelo.

Questão 3: Impacto do desbalanceamento no modelo

Se o dataset Titanic não estiver balanceado, o impacto será:

1. **Previsões enviesadas:** O modelo pode aprender a prever principalmente a classe majoritária (por exemplo, "não sobreviveu"), ignorando os padrões da classe minoritária ("sobreviveu").
2. **Métricas distorcidas:** A acurácia pode parecer alta, mas métricas como **recall** e **F1-score** para a classe minoritária seriam baixas.
3. **Soluções para lidar com o desbalanceamento:**
 - **Reamostragem:** Balancear o dataset usando técnicas como *oversampling* (aumentar a classe minoritária) ou *undersampling* (reduzir a classe majoritária).
 - **Atribuir pesos às classes:** Ajustar a função de perda para penalizar erros na classe minoritária.
 - **Estratégias de dados sintéticos:** Usar métodos como **SMOTE** (**S**ynthetic **M**inority **O**versampling **T**echnique) para gerar novos exemplos da classe minoritária.

GRUPO 4
(6 valores)

Comente as afirmações seguintes, assinalando a sua veracidade (**V**) ou falsidade (**F**), justificando a resposta EXCLUSIVAMENTE no espaço disponibilizado.

NÃO SÃO CONSIDERADAS respostas para as quais não exista justificação expressa.

QUESTÃO 1

- ☒ No desenvolvimento de sistemas de aprendizagem automática, a fase de preparação de dados tem particular importância porque os dados obtidos do «mundo físico» são incompletos, contêm lixo e são falsos.

Verdadeiro (V)

- **Justificativa:** Os dados obtidos do mundo físico frequentemente apresentam problemas como valores incompletos, ruídos e inconsistências. A preparação de dados é essencial para corrigir essas questões e garantir a qualidade do modelo.

QUESTÃO 2

- ☐ Técnicas de aprendizagem automática baseadas no desenvolvimento de árvores de decisão são utilizadas exclusivamente para a resolução de problemas de classificação.

Falso (F)

- **Justificativa:** Técnicas baseadas em árvores de decisão, como *Decision Tree* ou *Random Forest*, podem ser usadas tanto para **classificação** quanto para **regressão**, dependendo do problema.

QUESTÃO 3

- ☐ Paradigmas de aprendizagem com supervisão exigem maior intervenção humana do que qualquer outro paradigma uma vez que necessitam de quem desempenhe o papel de supervisor.

Verdadeiro (V)

- **Justificativa:** Nos paradigmas de aprendizagem supervisionada, é necessário ter dados rotulados, o que exige intervenção humana para classificar e organizar os dados previamente.

QUESTÃO 4

- ☒ O tratamento de valores nulos (*missing values*) existentes num *dataset* pode envolver a remoção de observações/registos ou de atributos/características.

Verdadeiro (V)

- **Justificativa:** Tratamento de valores nulos pode incluir tanto a remoção de registos que contenham valores ausentes quanto a exclusão de atributos que possuem muitos dados faltantes, dependendo da relevância para o modelo.

QUESTÃO 5

- ☐ A matriz de confusão à direita apresenta um valor de *accuracy* de $\frac{165}{150}$.

Falso (F)

- **Justificativa:** A matriz de confusão não indica que o valor de *accuracy* seja $\frac{165}{150}$. Para calcular a *accuracy*, usa-se a fórmula:

$$\text{Accuracy} = \frac{\text{predições corretas}}{\text{total de predições}} = \frac{50 + 100}{165} = 0,91 \text{ (91\%)}$$

n=165

		PREVISÃO		
		NÃO	SIM	
ATUAL	NÃO	50	10	60
	SIM	5	100	105
		55	110	

QUESTÃO 6

- ☐ Num processo de aprendizagem automática, a qualidade dos dados não afeta os resultados do processo uma vez que na fase de preparação de dados serão resolvidos todos os problemas como, por exemplo, ruído, *outliers*, dados falsos ou dados duplicados.

Falso (F)

- **Justificativa:** Embora a fase de preparação de dados seja fundamental, ela não pode resolver completamente problemas como ruído, outliers ou dados falsos, especialmente quando eles são persistentes ou complexos. A qualidade dos dados impacta diretamente o desempenho do modelo.