

Trainee Data Science - Lapisco

Módulo 1 - Atividade 1



SQL

Para resolver os exercícios, utilize a descrição do modelo relacional abaixo e os comandos SQL aprendidos. Se houver a necessidade de utilizar mais de uma tabela em uma consulta, utilize o INNER JOIN, em vez do produto cartesiano.

Questão 1

Insira 3 registros para cada tabela criada.

Solução

```
1 INSERT INTO tbAtor
2 VALUES
3 (1,"Matheus Rocha",20,"Fortaleza",2000.0,"M"),
4 (2,"Roberto",40,"Fortaleza",15000.0,"M"),
5 (3,"Iagson",30,"Fortaleza",10000.0,"M");
6
7 INSERT INTO tbNovela
8 VALUES
9 (1,"Avenida Brasil",'2002-08-16',"2100-08-16",8),
10 (2,"Chocolate Com Pimenta",'2002-08-16',"2100-08-16",10),
11 (3,"O Rei do Gado",'2002-08-16',"2100-08-16",7);
12
13
14 INSERT INTO tbCapitulo
15 VALUES
16 (1,"Capitulo 1","2010-08-16",1),
17 (2,"Capitulo 1","2010-08-16",2),
18 (3,"Capitulo 1","2010-08-16",3);
19
20 INSERT INTO tbPersonagem
21 VALUES
22 (1,"Personagem A",15,"Rico",1),
23 (2,"Personagem B",20,"Rico Demais",2),
24 (3,"Personagem C",30,"Pobre",3);
25
26 INSERT INTO tbNovelaPersonagem
27 VALUES
28 (1,1),
29 (2,2),
30 (3,3);
```

Questão 2

Encontre todas as novelas que tenham o valor do horário de exibição vazio.

Solução

```
1 SELECT *
2 FROM tbNovela t
3 WHERE t.horario_exibicao IS NULL;
```

Questão 3

Selecione o nome de todos os atores que morem em cidades que comecem com a letra “M”.

Solução

```
1 SELECT *
2 FROM tbAtor t
3 WHERE t.cidade_ator LIKE "M%";
```

Questão 4

Selecione todos os campos da tabela tbPersonagem ordenados por nome em ordem crescente.

Solução

```
1 SELECT *
2 FROM tbPersonagem t
3 ORDER BY t.nome_personagem ASC;
```

Questão 5

Selecione quantos capítulos existem por novela, retorne o nome da novela e a quantidade de capítulos para a novela.

Solução

```
1 SELECT n.nome_novela, COUNT(*) as qtd_capitulos
2 FROM tbNovela n
3 INNER JOIN tbCapitulo c
4 ON n.codigo_novela = c.codigo_novela
5 GROUP BY c.codigo_novela;
```

Questão 6

Encontre o nome de todas as novelas que tem mais de 40 capítulos.

Solução

```
1 SELECT n.nome_novela
2 FROM tbNovela n
3 INNER JOIN tbCapitulo c
4 ON n.codigo_novela = c.codigo_novela
5 GROUP BY c.codigo_novela
6 HAVING COUNT(*) > 40;
```

Pandas

Nas questões seguintes, utilizar as seguintes bases de dados: titanic, iris.

Questão 1

Verifique a presença de outliers (e.g., valor maior que 3 vezes o desvio padrão da base) na base de dados iris ou outra base de dados. Caso não exista, insira uma ou mais linhas e selecione elas utilizando a regra mencionada.

Solução

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import warnings
6 from IPython.display import display
7 warnings.filterwarnings("ignore")
8
9 iris = sns.load_dataset('iris')
10 iris.species = iris.species.astype('category')
11 iris.head()
12
13 cols = ['sepal_length', 'sepal_width', 'petal_length', 'petal_width']
14 for index, col in enumerate(cols):
15     print(f"--- Outliers para a coluna: {col} ---")
16
17     upper_limit = iris[col].mean() + np.abs(3 * iris[col].std())
18     lower_limit = iris[col].mean() - np.abs(3 * iris[col].std())
19
20     query = f"{col} > {upper_limit:.2f} | {col} < {lower_limit:.2f}"
21
22     outliers = iris.query(query)
23
24     # se nao encontrei outliers, crio eles
25     if outliers.empty:
26         print("\nOutliers nao encontrados, gerando outliers\n")
27
28         #criando outlier
29         new_row = {}
30         for new_col in cols:
```

```

31         new_row[new_col] = 0
32
33         new_row[col] = upper_limit + 1.0
34         new_row['species'] = 'setosa'
35
36         #inserindo outlier no dataframe
37         iris = iris.append(new_row, ignore_index = True)
38
39         outliers = iris.query(query)
40
41         display(outliers)
42         print("\n\n\n")

```

Questão 2

Faça um merge de dois dataframes e crie uma nova coluna para indicar a que dataframe pertence cada linha.

Solução

```

1  import pandas as pd
2  import numpy as np
3  import seaborn as sns
4  import matplotlib.pyplot as plt
5  import warnings
6  from IPython.display import display
7  warnings.filterwarnings("ignore")
8
9  df1 = pd.DataFrame({
10     'a': [1,2,3],
11     'b': [10,11,12]
12 })
13
14  df2 = pd.DataFrame({
15     'a': [4,5,6],
16     'b': [13,14,15]
17 })
18
19  df1['data_from'] = 'dataframe 1'
20  df2['data_from'] = 'dataframe 2'
21
22
23  df3 = pd.concat([df1,df2], ignore_index=True)
24  df3.species = df3.data_from.astype('category')
25
26  display(df3)

```

Resultado

É gerado um dataframe com os dados de ambos dataframes anteriores e também com um campo adicional, contendo a informação de onde tal informação veio:

Tabela 1: Dataset gerado através da concatenação dos datasets 1 e 2

	a	b	data.from
0	1	10	dataframe 1
1	2	11	dataframe 1
2	3	12	dataframe 1
3	4	13	dataframe 2
4	5	14	dataframe 2
5	6	15	dataframe 2

Questão 3

Selecione e mostre a quantidade de passageiros do titanic que possuem mais de 27 anos e que sobreviveram ao acidente.

Solução

```

1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import warnings
6 from IPython.display import display
7 warnings.filterwarnings("ignore")
8
9 titanic = sns.load_dataset('titanic')
10 titanic['survived'] = titanic['survived'].astype('boolean')
11 display(titanic.head())
12
13 q3 = titanic.query('age > 27 and survived == True')
14 display(q3.head())

```

Resultado

Tabela 2: 5 primeiros passageiros que possuem mais de 27 anos e sobreviveram ao acidente

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone
1	True	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False
3	True	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False
11	True	1	female	58.0	0	0	26.5500	S	First	woman	False	C	Southampton	yes	True
15	True	2	female	55.0	0	0	16.0000	S	Second	woman	False	NaN	Southampton	yes	True
21	True	2	male	34.0	0	0	13.0000	S	Second	man	True	D	Southampton	yes	True

Questão 4

Transforme os dados contínuos de idade do titanic em dados categóricos de acordo com a seguinte regra:

- 0 a 18 anos → Criança

- 18 a 65 anos → Adulto
- maior de 65 → Idoso

Solução

```

1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import warnings
6 from IPython.display import display
7 warnings.filterwarnings("ignore")
8
9 titanic = sns.load_dataset('titanic')
10 titanic['survived'] = titanic['survived'].astype('boolean')
11 display(titanic.head())
12
13 def person_type(age: int) -> str:
14     if age < 18:
15         return 'Child'
16     if age < 65:
17         return 'Adult'
18     return 'Old'
19
20 titanic['person_type'] = [person_type(age) for age in titanic.age]
21 sample = titanic.sample(n=10, random_state=1)
22
23 display(sample)

```

Resultado

Tabela 3: Amostra de 10 passageiros após classificação baseada em idade

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone	person_type
862	True	1	female	48.0	0	0	25.9292	S	First	woman	False	D	Southampton	yes	True	Adult
223	False	3	male	NaN	0	0	7.8958	S	Third	man	True	NaN	Southampton	no	True	Old
84	True	2	female	17.0	0	0	10.5000	S	Second	woman	False	NaN	Southampton	yes	True	Child
680	False	3	female	NaN	0	0	8.1375	Q	Third	woman	False	NaN	Queenstown	no	True	Old
535	True	2	female	7.0	0	2	26.2500	S	Second	child	False	NaN	Southampton	yes	False	Child
623	False	3	male	21.0	0	0	7.8542	S	Third	man	True	NaN	Southampton	no	True	Adult
148	False	2	male	36.5	0	2	26.0000	S	Second	man	True	F	Southampton	no	False	Adult
3	True	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False	Adult
34	False	1	male	28.0	1	0	82.1708	C	First	man	True	NaN	Cherbourg	no	False	Adult
241	True	3	female	NaN	1	0	15.5000	Q	Third	woman	False	NaN	Queenstown	yes	False	Old

Questão 5

Aplique o hold out utilizando diferentes porcentagens de divisão para separar a base iris em treino e teste (e.g. 60/40, 70/30, 80/20, 90/10). Crie vetores para armazenar os dados da divisão de treino e teste da seguinte forma:

- X_{train} → armazena os dados de treino
- X_{test} → armazena os dados de teste

- y_{train} → armazena a label correspondente dos dados de treino
- y_{test} → armazena a label correspondente dos dados de teste

Solução

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import warnings
6 from IPython.display import display
7 warnings.filterwarnings("ignore")
8 from sklearn.model_selection import train_test_split
9
10 iris = sns.load_dataset('iris')
11 iris.species = iris.species.astype('category')
12 iris.info()
13
14 X = iris.copy()
15 X = X.drop(['species'], axis = 1)
16 y = iris.species.copy()
17
18 #usando 20% dos dados para teste e 80% para treino
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.2)
```
