



Atividade Técnica

Continuação...

→ Pergunta 2

Crie uma query, considerando o SGBD MySQL, para exibir todos os dados de uma tabela de Pontos de Venda (tabela origem PONTO_VENDA_UNIDADE) e restringir apenas os pontos de venda que possuem sell in maior que 20.000 (cam

A	B	C	D	E	F	G	H	I	J
ID_DATASET	ID_PONTO_VENDA	NOME_PONTO_VEN	PERFIL_PONTO_VE	DATA	ID_LINHA_PRODUT	NOME_LINHA_PROI	MARCA_LINHA_PR	TIPO_COLETA	VALOR
1021	115	INVOLVES	ATACADAO	01/09/2020	398	BISCOITOS SORTI	INVOLVES ZERO	DISPONIBILIDADE	SIM
1022	115	INVOLVES	ATACADAO	02/09/2020	407	MARGARINA	INVOLVES ZERO	DISPONIBILIDADE	SIM
1023	115	INVOLVES	ATACADAO	03/09/2020	408	MANTEIGA	INVOLVES ZERO	DISPONIBILIDADE	SIM

O dataset enviado anexo ao email tem o nome **dataset_teste_de.csv**. Entendo que se trata do mesmo citado acima na descrição da questão, **PONTO_VENDA_UNIDADE**. Conforme o print da imagem acima, não há o campo **SELLIN** na tabela, pode ser que seja a coluna **VALOR**, mas não tenho certeza.

O dataset enviado tem apenas 40 linhas/registros, sendo que na coluna VALOR há campos 'SIM/NÃO' do tipo texto (20 linhas) e campos do tipo inteiro (min 0, max 6). O máximo da tabela é bem inferior ao solicitado na questão (20.000); portanto, mesmo que eu limpe e deixe apenas os registros de número inteiro, a tabela retorna vazia.

Pergunta2
▶ RUN
📄 SAVE QUERY
⬇️ DOWNLOAD
👤 SHARE
✅ This query wi

```

1 SELECT *
2 FROM `involves-422612.involves_datasets.involves-dataset`
3 WHERE VALOR NOT IN ('SIM', 'NÃO') AND CAST(VALOR AS INT64) > 20000
4 ORDER BY NOME_PONTO_VENDA;
5

```

Query results
📄 SAVE RESULTS
📊 EXPLORE DATA

JOB INFORMATION
RESULTS
CHART
JSON
EXECUTION DETAILS
EXECUTION GRAPH

There is no data to display.

🔍 Inserir_Visitas

▶ RUN

💾 SAVE QUERY ▾

⬇️ DOWNLOAD

👤 SHARE ▾

🕒 SCHEDULE

⚙️ MORE ▾

```
1 --Criar tabela
2 CREATE TABLE `involves-422612.involves_datasets.involves-dataset-visitas` (
3   ID_VISITA INT64 NOT NULL,
4   FK_PDV INT64 NOT NULL,
5   FL_VISITADO INT64 NOT NULL,
6   DATA_VISITA DATE NOT NULL
7 );
8
9 -- Se necessario limpar a tabela
10 TRUNCATE TABLE `involves-422612.involves_datasets.involves-dataset-visitas`
11
12 -- Inserir dados na tablea
13 INSERT INTO
14   `involves-422612.involves_datasets.involves-dataset-visitas` (ID_VISITA,
15     FK_PDV,
16     FL_VISITADO,
17     DATA_VISITA)
18 SELECT
19   id AS ID_VISITA,
20   CAST(115 + FLOOR(RAND() * 2) AS INT64) AS FK_PDV,
21   CAST(FLOOR(RAND() * 2) AS INT64) AS FL_VISITADO,
22   DATE_ADD('2024-01-01', INTERVAL CAST(FLOOR(RAND() * 365) AS INT64) DAY) AS DATA_VISITA,
23 FROM
24   UNNEST(GENERATE_ARRAY(1, 100)) AS id;
```

Resposta:

🔍 Pergunta4

▶ RUN

💾 SAVE QUERY ▾

⬇️ DOWNLOAD

👤 SHARE ▾

🕒 SCHEDULE

⚙️ MORE ▾

✅ This que...

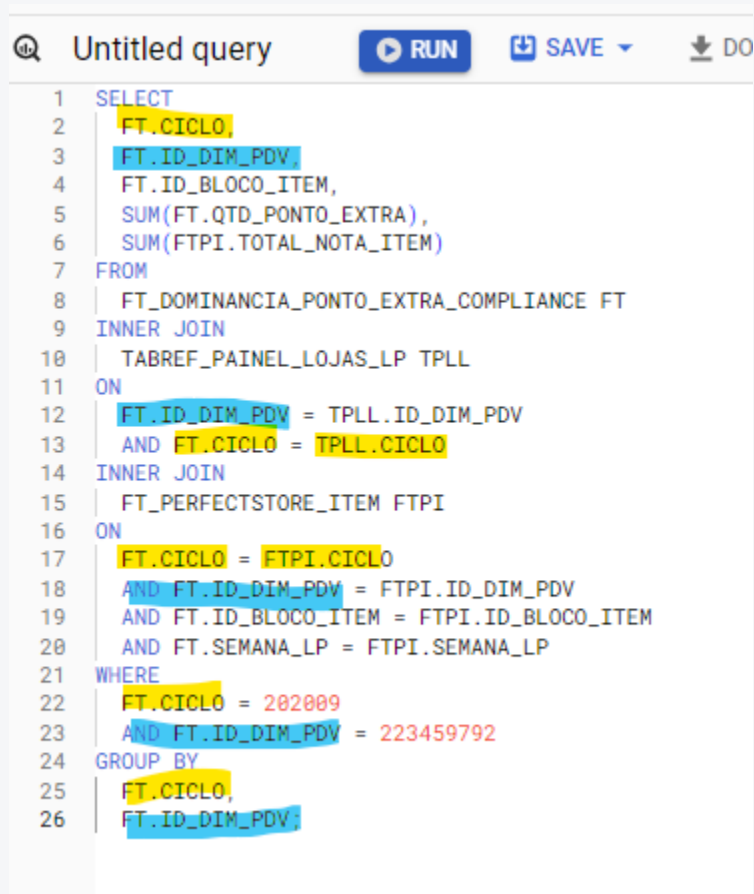
```
1 SELECT
2   PDV_UNIDADES.NOME_PONTO_VENDA,
3   COALESCE(VISITAS_COUNT.TOTAL_VISITAS, 0) AS TOTAL_VISITAS
4 FROM
5   `involves-422612.involves_datasets.involves-dataset` AS PDV_UNIDADES
6 LEFT JOIN (
7   SELECT
8     FK_PDV,
9     COUNT(ID_VISITA) AS TOTAL_VISITAS
10  FROM
11    `involves-422612.involves_datasets.involves-dataset-visitas`
12  WHERE
13    FL_VISITADO = 1
14  GROUP BY
15    FK_PDV) AS VISITAS_COUNT
16 ON
17   PDV_UNIDADES.ID_PONTO_VENDA = VISITAS_COUNT.FK_PDV
18 WHERE
19   NOME_PONTO_VENDA = 'INVOLVES'
20 GROUP BY
21   PDV_UNIDADES.NOME_PONTO_VENDA,
22   TOTAL_VISITAS;
```

Query results		SAVE RESULTS ▾	EXPLORE DATA ▾	↺
JOB INFORMATION		RESULTS	CHART	JSON
		EXECUTION DETAILS	EXECUTION GRAPH	
Row	NOME_PONTO_VENDA ▾	TOTAL_VISITAS ▾		
1	INVOLVES	22		

→ Pergunta 5

Considerando a query abaixo, a pessoa engenheira de dados identificou que a performance da query está muito abaixo do esperado. Imaginando que um dos problemas possa estar relacionado aos índices das tabelas do banco de dados, a pessoa resolveu analisar os índices nas tabelas. Liste quais possíveis campos devem ser indexados nas tabelas do banco de dados para que a query tenha a melhor performance possível. Leve em consideração que o primeiro campo no banco de dados está indexado.

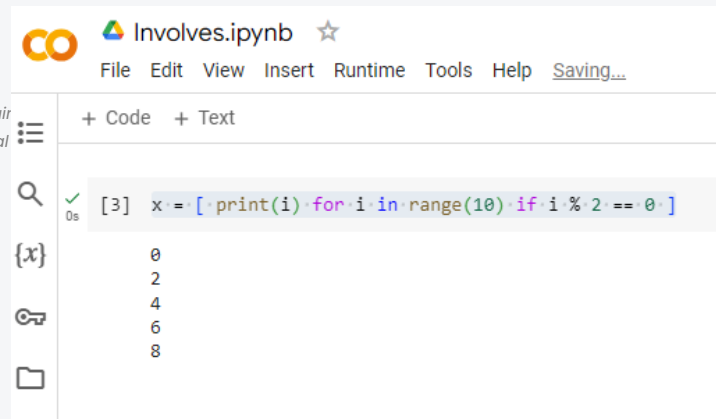
A indexação agiliza o encontro entre dois campos e evita a necessidade de uma consulta ser realizada linha por linha. Vários campos podem ser indexados, mas a prioridade na consulta apresentada deve ser para os campos que estão nas junções, nos filtros e nos agrupamentos. Abaixo, podemos observar que os campos FT.CICLO e FT.ID_DIM_PDV estão presentes em todas as junções e nos filtros. Portanto, são eles que precisam ser indexados com urgência, inclusive com a possibilidade de ser criado um índice composto.



```
1 SELECT
2   FT.CICLO,
3   FT.ID_DIM_PDV,
4   FT.ID_BLOCO_ITEM,
5   SUM(FT.QTD_PONTO_EXTRA),
6   SUM(FTPI.TOTAL_NOTA_ITEM)
7 FROM
8   FT_DOMINANCIA_PONTO_EXTRA_COMPLIANCE FT
9 INNER JOIN
10  TABREF_PAINEL_LOJAS_LP TPLL
11 ON
12  FT.ID_DIM_PDV = TPLL.ID_DIM_PDV
13  AND FT.CICLO = TPLL.CICLO
14 INNER JOIN
15  FT_PERFECTSTORE_ITEM FTPI
16 ON
17  FT.CICLO = FTPI.CICLO
18  AND FT.ID_DIM_PDV = FTPI.ID_DIM_PDV
19  AND FT.ID_BLOCO_ITEM = FTPI.ID_BLOCO_ITEM
20  AND FT.SEMANA_LP = FTPI.SEMANA_LP
21 WHERE
22  FT.CICLO = 202009
23  AND FT.ID_DIM_PDV = 223459792
24 GROUP BY
25  FT.CICLO,
26  FT.ID_DIM_PDV;
```

→ Pergunta 6

Considere a instrução Python a seguir no Python, a variável "x" conterá qual



The screenshot shows a Jupyter Notebook titled "Involves.ipynb". The code cell contains a list comprehension: `x = [print(i) for i in range(10) if i % 2 == 0]`. The output cell shows the result of the list comprehension, which is an empty list `[]`, and the printed values of `i` (0, 2, 4, 6, 8) on separate lines.

```
[3] x = [print(i) for i in range(10) if i % 2 == 0.]
```

```
0
2
4
6
8
[]
```

→ Pergunta 7

Faça um script em Python que peça dois números e imprima a soma

```
[11] # @title Default title text
# Apresenta
print ('Olá, vamos fazer a soma de dois números')

# Pede um 1º número
num1 = float(input("Digite o 1º número: "))

# Pede um 2º número
num2 = float(input("Digite o 2º número: "))

# Calcula a soma dos dois números inseridos
soma = num1 + num2

# Imprime o resultado
print("O resultado da soma é:", soma)
```

```
Olá, vamos fazer a soma de dois números
Digite o 1º número: 50
Digite o 2º número: 150
O resultado da soma é: 200.0
```

Para responder às questões 8, 9 e 10

A ETL final deve conter um job princ.

10. Além disso, que tal ganhar um ponto a mais nessas questões? Para isso, inclua o projeto criado em um repositório do Github (é importante que seja público para termos visibilidade, ok?). Compartilhe por aqui o link para o repositório.

→ Pergunta 8

Construa uma transformação que deve usar como datasource o dataset (DATASET_TESTE_DE.csv) que contém informações de coletas de dados nos ponto de vendas. A ETL deve consultar o dataset e inserir, em uma base de dados (modo de simulação), as informações das coletas no formato abaixo:

- A pergunta a) não especifica o formato do ano e mês (ex. YY, YYYY-MM, MMM, etc..), deixei o padrão da função. Não encontrei o campo perfil da linha de produto, adicionei MARCA_LINHA_PRODUTO
- a) Dimensão Calendário (DIM_CALENDARIO): Deve conter data, mês e ano da coleta
 - b) Dimensão Ponto de Venda (DIM_PONTO_VENDA): Deve conter o ID, nome e perfil do ponto de venda
 - c) Dimensão Linha de Produto (DIM_LINHA_PRODUTO): Deve conter o id, nome e perfil da linha de produto

Pergunta8_DIM_CALENDARIO

```
1 SELECT
2   DATA,
3   EXTRACT(YEAR
4     FROM
5       DATA) AS ANO,
6   EXTRACT(MONTH
7     FROM
8       DATA) AS MES
9 FROM
10  `involves-422612.involves_datasets.involves-dataset`
```

Pergunta8_DIM_LINHA_PRODU...

```
1 SELECT
2   ID_LINHA_PRODUTO,
3   NOME_LINHA_PRODUTO,
4   MARCA_LINHA_PRODUTO
5 FROM
6  `involves-422612.involves_datasets.involves-dataset`
```

Pergunta8_DIM_PDV

```
1 SELECT
2   ID_PONTO_VENDA,
3   NOME_PONTO_VENDA,
4   PERFIL_PONTO_VENDA
5 FROM
6  `involves-422612.involves_datasets.involves-dataset`
```

→ Pergunta 9

Construa uma transformação que deve usar como datasource o dataset (DATASET_TESTE_DE.csv) que contém informações de coletas de dados nos ponto de vendas. A transformação deve consultar o dataset e inserir, em uma base de dados (modelo dimensional), as informações coletadas, conforme as tabelas abaixo:

a) Fato Disponibilidade (FT_DISPONIBILIDADE): Deve conter os ids de ligação das tabelas de dimensões criadas na questão anterior e a quantidade de presenças de cada linha de produto no mês de Setembro/20.

b) Fato Disponibilidade Agregada (FT_DISPONIBILIDADE_AGREGADA): Deve conter os ids de ligação das tabelas de dimensões agregadas criadas na questão anterior e a quantidade de presenças de cada linha de produto por ponto de venda no mês de Setembro/20.

Não agreguei por DATA, porque conforme descrição apenas os dados de setembro são neccessarios.

```

Pergunta9_FT_DISPONIBILIDADE
RUN

1 SELECT
2   ID_PONTO_VENDA,
3   ID_LINHA_PRODUTO,
4   COUNT(*) AS QTD_PRESENCIA
5 FROM
6   'involves-422612.involves_datasets.involves-dataset'
7 WHERE
8   EXTRACT(YEAR
9     FROM
10    | DATA) = 2020
11   AND EXTRACT(MONTH
12     FROM
13    | DATA) = 9
14   AND VALOR = 'SIM'
15 GROUP BY
16   ID_PONTO_VENDA,
17   ID_LINHA_PRODUTO
18 ORDER BY
19   ID_LINHA_PRODUTO,
20   ID_PONTO_VENDA,
21   QTD_PRESENCIA DESC;
```

Query results

JOB INFORMATION		RESULTS	CHART	JSON
Row	ID_PONTO_VENDA	ID_LINHA_PRODUTO	QTD_PRESENCIA	
1	115	398	2	
2	115	407	1	
3	116	407	2	
4	115	408	1	
5	115	422	1	
6	116	422	2	
7	115	423	2	

Pergunta9_FT_DISPONIBILIDADE_AGREGA...

```

1 SELECT
2   ID_PONTO_VENDA,
3   COUNT(*) AS QTD_PRESENCA
4 FROM
5   `involves-422612.involves_datasets.involves-dataset`
6 WHERE
7   EXTRACT(YEAR
8     FROM
9       | DATA) = 2020
10  AND EXTRACT(MONTH
11    FROM
12      | DATA) = 9
13  AND VALOR = 'SIM'
14 GROUP BY
15   ID_PONTO_VENDA
16 ORDER BY
17   ID_PONTO_VENDA,
18   QTD_PRESENCA DESC;
19
    
```

Query results

JOB INFORMATION		RESULTS	CHART
Row	ID_PONTO_VENDA	QTD_PRESENCA	
1	115	7	
2	116	4	