


```
[70]: #-----CÓDIGO COMENTADO-----
# Neste caso vamos usar a operação drill-across que compara medidas numéricas de tabelas de fatos de
# diferentes
# utilizando uma dimensão em comum (no caso data). Neste caso, vamos precisar separar as consultas e
# em dois blocos,
# sendo o primeiro referente aos dados de salário e o segundo referente aos dados de receitas.

# 01. utilizando dados de pagamento e atribuindo ao bloco 1 nomeado de pag
#
# 02. juntando com dados de data pela chave primária dataPK
#
# 03. juntando com dados de funcionario pela chave primária funcPK
#
# 04. aplicando condição para identificar funcionários do estado do RIO DE JANEIRO (tanto por nome quan
to por sigla)
#
# 05. selecionando os dados de dataAno e salario como resultados do bloco 1
#
# 06. agrupando os dados por dataAno
#
# 07. agregando os dados de salário como o total do salário das funcionárias
#
# 08. utilizando dados de negociacao e atribuindo ao bloco 2 nomeado de neg
#
# 09. juntando com dados de data pela chave primária dataPK
#
# 10. juntando com dados de equipe pela chave primária equipePK
#
# 11. aplicando condição para identificar equipes do estado do RIO DE JANEIRO (tanto por nome quanto po
r sigla)
#
# 12. selecionando os dados de dataAno e receita como resultados do bloco 2
#
# 13. agrupando os dados por dataAno
#
# 14. agregando os dados de receita como o total da receita das equipes
#
# 15. utilizando dados do bloco 1 nomeado de pag
#
# 16. juntando com dados do bloco 2 nomeado de neg pela dimensão em comum entre ambas, no caso, dataAno
#
# 17. selecionando os dados de dataAno, total de salários e total de receitas
#
# 18. arredondando a soma dos salários para 2 casas decimais
#
# 19. arredondando a soma das receitas para 2 casas decimais
#
# 20. renomeando dataAno para ANO
#
# 21. renomeando sum(salario) para TOTALSALARIO
#
# 22. renomeando sum(receita) para TOTALRECEITA
#
# 23. ordenando por data, no caso ANO de forma ascendente
#
# 24. apresentando os 20 primeiros resultados sem truncamento das strings
#
#-----#
pag = pagamento\
    .join(data, on='dataPK')\
    .join(funcionario, on='funcPK')\
    .where('funcSexo = "F" AND (funcEstadoNome = "RIO DE JANEIRO" OR funcEstadoSigla = "RJ")')\
    .select('dataAno', 'salario')\
    .groupBy('dataAno')\
    .sum('salario')

neg = negociacao\
    .join(data, on='dataPK')\
    .join(equipe, on='equipePK')\
    .where('filialEstadoNome = "RIO DE JANEIRO" OR filialEstadoSigla = "RJ"')\
    .select('dataAno', 'receita')\
    .groupBy('dataAno')\
    .sum('receita')

pag\
    .join(neg, on='dataAno')\
    .select('dataAno', 'sum(salario)', 'sum(receita)')\
    .withColumn('sum(salario)', round('sum(salario)',2))\
    .withColumn('sum(receita)', round('sum(receita)',2))\
    .withColumnRenamed('dataAno', 'ANO')\
    .withColumnRenamed('sum(salario)', 'TOTALSALARIO')\
    .withColumnRenamed('sum(receita)', 'TOTALRECEITA')\
    .orderBy('ANO',ascending=True)\
    .show(20,truncate=False)

#-----#
ANO TOTALSALARIO TOTALRECEITA
-----#
(2016|30061.2 | 12205042.91 |
(2017|30061.2 | 13484981.8 |
(2018|48108.36 | 14741199.75 |
(2019|70794.36 | 15100984.61 |
(2020|70794.36 | 14192420.2 |
-----#
```

O objetivo da análise desta seção é obter uma visão relacionada aos sexos, por meio da comparação do total anual de gastos em salários para o pagamento das mulheres e dos homens em comparação ao total anual de receitas recebidas.

Liste, para cada dataAno, a soma dos salários das funcionárias de sexo feminino, a soma dos salários dos funcionários do sexo masculino e a soma das receitas recebidas. Arredonde a soma dos salários e a soma das receitas para até duas casas decimais. Devem ser exibidas as colunas na ordem e com os nomes especificados a seguir: "ANO", "TOTALSALARIO/MULHERES", "TOTALSALARIO/HOMENS", "TOTALRECEITA". Ordene as linhas exibidas por ano em ordem ascendente. Liste as primeiras 20 linhas da resposta, sem truncamento das strings.

Consulta OLAP na linguagem SQL.

```
In [78]: #-----CÓDIGO COMENTADO-----
# Neste caso vamos usar a operação drill-across que compara medidas numéricas de tabelas de fatos di
ferentes
# utilizando uma dimensão em comum (no caso data).

# 01. selecionando dataAnoFF (dataAno já consolidada na subquery de salário para sexo feminino) como AN
O
#
# selecionando a soma dos salários femininos arredondando com 2 casas como TOTALSALARIO/MULHERES
#
# selecionando a soma dos salários masculinos arredondando com 2 casas como TOTALSALARIO/HOMENS
#
# selecionando a soma das receitas arredondando com 2 casas como TOTALRECEITA
#
# Nestes casos, tivemos alguns números sendo apresentados com notação científica, portanto, aplicamo
s a função CAST
# AS DECIMAL, para a saída ficar no formato adequado (requerido no exercício). Como parâmetros do DE
CIMAL passamos o
# valor 10 (tamanho considerado adequado - número de dígitos permitidos na parte inteira) e 2 (númer
o de casas decimais).
# Vamos executar três sub-selects no comando FROM. O primeiro responsável pelos dados de salário fem
inino
# o segundo responsável pelos dados de salário masculino e o terceiro responsável pelos dados de rec
eita.
# 02. selecionando dataAno de data e soma dos salários como salarioF (que terá os dados para funcioná
rias de sexo feminino)
#
# 03. especificando as relações de pagamento com data contendo mesma dataPK
#
# 04. especificando as relações de funcionario com pagamento contendo mesmo funcPK
#
# 05. aplicando condição para funcionárias apenas do sexo feminino
#
# 06. agrupando os dados por dataAno
#
# 07. aplicando os resultados do sub-select como dataAnoFF para dataAno e salarioF como o total do sala
rio das mulheres
#
# 08. cláusula para juntarmos o segundo sub-select
#
# 09. selecionando dataAno de data e soma dos salários como salarioM (que terá os dados para funcioná
rias de sexo masculino)
#
# 10. especificando as relações de pagamento com data contendo mesma dataPK
#
# 11. especificando as relações de funcionario com pagamento contendo mesmo funcPK
#
# 12. aplicando condição para funcionários apenas do sexo masculino
#
# 13. agrupando os dados por dataAno
#
# 14. aplicando os resultados do sub-select como dataAnoPM para dataAno e salarioM como o total do sala
rio dos homens
#
# 15. cláusula para juntarmos o terceiro sub-select
#
# 16. selecionando dataAno de data e soma das receitas de negociação
#
# 17. especificando as relações de data com negociação contendo mesma dataPK
#
# 18. especificando as relações de equipe com negociação contendo mesmo equipePK
#
# 19. agrupando os dados por dataAno
#
# 20. aplicando os resultados do sub-select como dataAnoN para dataAno e receita como o total do receit
a das equipes
#
# 21. aplicando condição para juntar a dimensão em comum, dataAnoFF igual ao dataAnoN, igual ao dataAn
oPM
#
# 22. agrupando por data, no caso dataAnoFF
#
# 23. ordenando por data, no caso dataAnoFF de forma ascendente
#
# 24. apresentando os 20 primeiros resultados sem truncamento das strings
#
#-----#

consultaSQL = """
SELECT sum(dataAnoFF AS 'ANO', CAST(ROUND(SUM(salarioF),2) AS DECIMAL(10,2)) AS 'TOTALSALARIO/MULHERES', CAS
T(ROUND(SUM(salarioM),2) AS DECIMAL(10,2)) AS 'TOTALSALARIO/HOMENS', CAST(ROUND(SUM(receita),2) AS DECIM
AL(10,2)) AS 'TOTALRECEITA
FROM (
    SELECT dataAno, SUM(salario) AS 'salarioF'
    FROM pagamento JOIN data ON (pagamento.dataPK = data.dataPK)
    JOIN funcionario ON (funcionario.funcPK = pagamento.funcPK)
    WHERE funcSexo = 'F'
    GROUP BY dataAno
) AS salF(dataAnoFF, salarioF)
JOIN (
    SELECT dataAno, SUM(salario) AS 'salarioM'
    FROM pagamento JOIN data ON (pagamento.dataPK = data.dataPK)
    JOIN funcionario ON (funcionario.funcPK = pagamento.funcPK)
    WHERE funcSexo = 'M'
    GROUP BY dataAno
) AS salM (dataAnoPM, salarioM)
JOIN (
    SELECT dataAno, SUM(receita)
    FROM data JOIN negociacao ON (data.dataPK = negociacao.dataPK)
    JOIN equipe ON (equipe.equipePK = negociacao.equipePK)
    GROUP BY dataAno
) AS rec(dataAnoN, receita)
WHERE (dataAnoFF = dataAnoN AND dataAnoFF = dataAnoPM)
GROUP BY dataAnoFF
ORDER BY dataAnoFF ASC
"""

spark.sql(consultaSQL).show(20,truncate=False)

#-----#
ANO TOTALSALARIO/MULHERES TOTALSALARIO/HOMENS TOTALRECEITA
-----#
(2016|1210393.21 | 3232223.89 | 4614246.97 |
(2017|2500857.87 | 7204221.87 | 7200423.35 |
(2018|3800427.49 | 11135098.98 | 11153333.66 |
(2019|4733247.25 | 13834419.20 | 35353318.33 |
(2020|4733247.25 | 13834419.20 | 30222175.87 |
-----#
```