





```
[43]: # Resposta da Questão 4

#-----CÓDIGO COMENTADO-----#
# 01. selecionando dataAno como ANO
#
# aplicando um select para extrair o ano atual e subtraindo de funcAnoNascimento para definir a IDA
# selecionando funcRegiaoNome como REGIAO
# realizando a soma dos salários arredondando para 2 casas decimais como TOTALSALARIO
# 02. especificando as relações de data com pagamento contendo mesma dataPK
#
# 03. especificando as relações de funcionario com pagamento contendo mesmo funcPK
#
# 04. aplicando condição para funcionárias apenas do sexo feminino
#
# 05. que nasceram entre 1970 (inclusive) e 1990 (inclusive) aplicando a função BETWEEN
#
# 06. que vivem no SUDESTE ou no NORDESTE (tanto pelo nome da região quanto pela sigla)
#
# 07. agrupando os dados por ANO, IDADE e REGIAO
#
# 08. ordenando por ANO, IDADE e REGIAO de forma ascendente
#
# 09. apresentando os 20 primeiros resultados sem truncamento das strings
#
#-----#

consultaSQL = """
SELECT dataAno AS 'ANO', ((SELECT EXTRACT (YEAR FROM CURRENT_DATE)) - funcAnoNascimento) AS 'IDADE', fu
ncRegiaoNome AS 'REGIAO', ROUND(SUM(salario),2) AS 'TOTALSALARIO'
FROM data JOIN pagamento ON (data.dataPK = pagamento.dataPK)
FROM funcionario ON (funcionario.funcPK = pagamento.funcPK)
WHERE funcSexo = 'F'
AND funcAnoNascimento BETWEEN 1970 AND 1990
AND funcRegiaoNome = 'SUDESTE' OR funcRegiaoNome = 'NORDESTE' OR funcRegiaoSigla = 'SE' OR func
RegiaoSigla = 'NE')
GROUP BY ANO, IDADE, REGIAO
ORDER BY ANO ASC, IDADE ASC, REGIAO ASC
"""
spark.sql(consultaSQL).show(20,truncate=False)
```

Questão 4

(valor:1.5) Considere que as equipes cujos valores de 'equipePK' são iguais a 1, 3 e 5 possuem a maior quantidade de funcionários do sexo feminino. Liste, para cada dataAno, a soma das receitas recebidas por essas equipes, o nome da equipe, o nome da filial e o setor do cliente, considerando apenas os clientes localizados na cidade de "SAO CARLOS". Arredonde a média dos salários para até duas casas decimais. Devem ser exibidas as colunas na ordem e com os nomes especificados a seguir: "ANO", "NOMEEQUIPE", "NOMEFILIAL", "SETORCLIENTE", "TOTALRECEITA". Ordene as linhas exibidas primeiro por ano, depois por nome da equipe, depois por nome da filial e depois por setor do cliente, todos em ordem ascendente. Liste as primeiras 20 linhas da resposta, sem truncamento das strings.

Resolva a questão especificando a consulta usando os métodos de pyspark.sql.

Resposta da Questão 4

```
In [42]: # Resposta da Questão 4

#-----CÓDIGO COMENTADO-----#
# 01. utilizando dados de negociacao
#
# 02. juntando com dados de equipe pela chave primária equipePK
#
# 03. juntando com dados de data pela chave primária dataPK
#
# 04. juntando com dados de cliente pela chave primária clientePK
#
# 05. aplicando condição para equipePK sendo 1, 3 e 5 e clientes da cidade de SAO CARLOS
#
# 06. selecionando os dados requisitados dataAno, equipeNome, filialNome, clienteSetor e receita
#
# 07. agrupando por dataAno, equipeNome, filialNome, clienteSetor e agregando pela soma das receitas
#
# 08. arredondando a soma da receitas para 2 casas decimais
#
# 09. renomeando dataAno para ANO
#
# 10. renomeando equipeNome para NOMEEQUIPE
#
# 11. renomeando filialNome para NOMEFILIAL
#
# 12. renomeando clienteSetor para SETORCLIENTE
#
# 13. renomeando sum(receita) para TOTALRECEITA
#
# 14. ordenando por ANO, NOMEEQUIPE, NOMEFILIAL e SETORCLIENTE de forma ascendente
#
# 15. apresentando os 20 primeiros resultados sem truncamento das strings
#
#-----#

negociacao\
.join(equipe, on='equipePK')\
.join(data, on='dataPK')\
.join(cliente, on='clientePK')\
.where('(equipePK = 1 OR equipePK = 3 OR equipePK = 5) AND clienteCidade = "SAO CARLOS"')\
.select('dataAno', 'equipeNome', 'filialNome', 'clienteSetor', 'receita')\
.groupBy('dataAno', 'equipeNome', 'filialNome', 'clienteSetor').sum('receita')\
.withColumn('sum(receita)',round('sum(receita)',2))\
.withColumnRenamed('dataAno','ANO')\
.withColumnRenamed('equipeNome','NOMEEQUIPE')\
.withColumnRenamed('filialNome','NOMEFILIAL')\
.withColumnRenamed('clienteSetor','SETORCLIENTE')\
.withColumnRenamed('sum(receita)','TOTALRECEITA')\
.orderBy('ANO','NOMEEQUIPE','NOMEFILIAL','SETORCLIENTE',ascending=True)\
.show(20,truncate=False)
```

6.3 Visão Geral da Atuação Feminina

O objetivo das análises desta seção é obter uma visão relacionada especificamente à atuação feminina, considerando aspectos conjuntos referentes a salários e receitas. Podem ser realizadas diferentes análises, dentre as quais destaca-se a análise base descrita a seguir.

Análise Base

Liste, para cada dataAno, a soma dos salários das funcionárias de sexo feminino que moram no estado do "RIO DE JANEIRO" ("RJ") e as somas das receitas recebidas pelas equipes localizadas no estado do "RIO DE JANEIRO" ("RJ"). O estado no qual as funcionárias moram pode ser identificado pelos atributos funcEstadoNome ou funcEstadoSigla da tabela de dimensão funcionario.

enquanto que o estado nos quais as equipes estão localizadas pode ser identificado pelos atributos filialEstadoNome ou filialEstadoSigla da tabela de dimensão equipe. 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", "TOTALRECEITA". Ordene as linhas exibidas por ano em ordem ascendente. Liste as primeiras 20 linhas da resposta, sem truncamento das strings.

Questão 5

(valor: 1.5) Resolva a "Análise Base" especificando a consulta OLAP na linguagem SQL.

Resposta da Questão 5

```
In [64]: # Resposta da Questão 5

#-----CÓDIGO COMENTADO-----#
# Neste caso vamos usar a operação drill-across que compara medidas numéricas de tabelas de fatos de diferentes tabelas
# utilizando uma dimensão em comum (no caso data).
#
# 01. selecionando dataAnoP (dataAno já consolidada na subquery de salário) como ANO
#
# selecionando a soma dos salários arredondando com 2 casas como TOTALSALARIO
#
# selecionando a soma das receitas arredondando com 2 casas como TOTALRECEITA
#
# Vamos executar dois sub-selects no comando FROM. O primeiro responsável pelos dados de salário e o segundo de receita
#
# 02. selecionando dataAnoP de data e soma dos salários de funcionario
#
# 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. aplicando condição para funcionárias do RIO DE JANEIRO (filtrando tanto pelo nome do estado quanto pela sigla)
#
# 07. agrupando os dados por dataAno
#
# 08. aplicando os resultados do sub-select como dataAnoP para dataAno e salário como o total do salário das funcionárias
#
# 09. cláusula para fazer a junção com outro sub-select, agora aplicado nos dados de receita
#
# 10. selecionando dataAnoP de data e soma das receitas de negociação
#
# 11. especificando as relações de data com negociação contendo mesma dataPK
#
# 12. especificando as relações de equipe com negociação contendo mesmo equipePK
#
# 13. aplicando condição para equipes do RIO DE JANEIRO (filtrando tanto pelo nome do estado quanto pela sigla)
#
# 14. agrupando os dados por dataAno
#
# 15. aplicando os resultados do sub-select como dataAnoP para dataAno e receita como o total do receita das equipes
#
# 16. aplicando condição para juntar a dimensão em comum, dataAnoP do salário igual ao dataAnoP da receita
#
# 17. agrupando por data, no caso dataAnoP
#
# 18. ordenando por data, no caso dataAnoP de forma ascendente
#
# 19. apresentando os 20 primeiros resultados sem truncamento das strings
#
#-----#

consultaSQL = """
SELECT dataAnoP AS 'ANO', ROUND(SUM(salario),2) AS 'TOTALSALARIO', ROUND(SUM(receita),2) AS 'TOTALRECEITA'
FROM ( SELECT dataAno, SUM(salario)
FROM pagamento JOIN data ON (pagamento.dataPK = data.dataPK)
FROM funcionario ON (funcionario.funcPK = pagamento.funcPK)
WHERE funcSexo = 'F'
AND funcEstadoNome = 'RIO DE JANEIRO' OR funcEstadoSigla = 'RJ')
GROUP BY dataAno
) AS sal(dataAnoP, salario)
JOIN
( SELECT dataAno, SUM(receita)
FROM data JOIN negociacao ON (data.dataPK = negociacao.dataPK)
FROM funcionario ON (equipe.equipePK = negociacao.equipePK)
WHERE filialEstadoNome = 'RIO DE JANEIRO' OR filialEstadoSigla = 'RJ')
GROUP BY dataAno
) AS rec(dataAnoN, receita)
WHERE dataAnoP = dataAnoN
GROUP BY dataAnoP
ORDER BY dataAnoP ASC
"""
spark.sql(consultaSQL).show(20,truncate=False)
```

ANO	TOTALSALARIO	TOTALRECEITA
2016	30061.2	12205042.91
2017	30061.2	13484981.8
2018	148108.36	14741199.75
2019	10794.36	15100984.61
2020	10794.36	14192420.2

Questão 6

(valor: 1.5) Resolva a "Análise Base" especificando a consulta OLAP usando os métodos de pyspark.sql.

Resposta da Questão 6

```
In [70]: # Resposta da Questão 6

#-----CÓDIGO COMENTADO-----#
# Neste caso vamos usar a operação drill-across que compara medidas numéricas de tabelas de fatos de diferentes tabelas
# utilizando uma dimensão em comum (no caso data). Neste caso, vamos precisar separar as consultas 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 funcionárias do estado do RIO DE JANEIRO (tanto por nome quanto 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 por 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')\
.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	148108.36	14741199.75
2019	10794.36	15100984.61
2020	10794.36	14192420.2

6.4 Visão Comparativa Final

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.

Questão 7

(valor:2.0) 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", "TOTALRECEITA", "TOTALSALARIOHOMENS", "TOTALRECEITA". Ordene as linhas exibidas por ano em ordem ascendente. Liste as primeiras 20 linhas da resposta, sem truncamento das strings.

Resolva a questão especificando a consulta OLAP na linguagem SQL.

Resposta da Questão 7

```
In [78]: # Resposta da Questão 7

#-----CÓDIGO COMENTADO-----#
# Neste caso vamos usar a operação drill-across que compara medidas numéricas de tabelas de fatos de diferentes tabelas
# utilizando uma dimensão em comum (no caso data).
#
# 01. selecionando dataAnoP (dataAno já consolidada na subquery de salário para sexo feminino) como ANO
#
# selecionando a soma dos salários femininos arredondando com 2 casas como TOTALSALARIOMULHERES
#
# selecionando a soma dos salários masculinos arredondando com 2 casas como TOTALSALARIOHOMENS
#
# selecionando a soma das receitas arredondando com 2 casas como TOTALRECEITA
#
# Nestes casos, tivemos alguns números sendo apresentados com notação científica, portanto, aplicamos a função CAST para a saída ficar no formato adequado (requerido no exercício). Como parâmetros do DE CIMAL passamos o # considerado adequado - número de dígitos permitidos na parte inteira) e 2 (número de casas decimais).
#
# Vamos executar três sub-selects no comando FROM. O primeiro responsável pelos dados de salário feminino, o segundo responsável pelos dados de salário masculino e o terceiro responsável pelos dados de receita.
#
# 02. selecionando dataAnoP de data e soma dos salários como salarioF (que terá os dados para funcionárias de sexo feminino)
#
# 03. selecionando 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 dataAnoPF para dataAno e salarioF como o total do salário das mulheres
#
# 08. cláusula para juntarmos o segundo sub-select
#
# 09. selecionando dataAnoP de data e soma dos salários como salarioM (que terá os dados para funcionários 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 salário dos homens
#
# 15. cláusula para juntarmos o terceiro sub-select
#
# 16. selecionando dataAnoP 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 receita das equipes
#
# 21. aplicando condição para juntar a dimensão em comum, dataAnoPF igual ao dataAnoN, igual ao dataAnoPM
#
# 22. agrupando por data, no caso dataAnoPF
#
# 23. ordenando por data, no caso dataAnoPF de forma ascendente
#
# 24. apresentando os 20 primeiros resultados sem truncamento das strings
#
#-----#

consultaSQL = """
SELECT dataAnoPF AS 'ANO', CAST(ROUND(SUM(salarioF),2) AS DECIMAL(10,2)) AS 'TOTALSALARIO', CAST(ROUND(SUM(salarioM),2) AS DECIMAL(10,2)) AS 'TOTALSALARIOHOMENS', CAST(ROUND(SUM(receita),2) AS DECIMAL(10,2)) AS 'TOTALRECEITA'
FROM ( SELECT dataAno, SUM(salarioF) AS 'salarioF'
FROM pagamento JOIN data ON (pagamento.dataPK = data.dataPK)
FROM funcionario ON (funcionario.funcPK = pagamento.funcPK)
WHERE funcSexo = 'F'
GROUP BY dataAno
) AS salF(dataAnoPF, salarioF)
JOIN
( SELECT dataAno, SUM(salarioM) AS 'salarioM'
FROM pagamento JOIN data ON (pagamento.dataPK = data.dataPK)
FROM 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)
FROM funcionario ON (equipe.equipePK = negociacao.equipePK)
GROUP BY dataAno
) AS rec(dataAnoN, receita)
WHERE dataAnoPF = dataAnoN AND dataAnoPM = dataAnoN
GROUP BY dataAnoPF
ORDER BY dataAnoPF ASC
"""
spark.sql(consultaSQL).show(20,truncate=False)
```

ANO	TOTALSALARIO	TOTALSALARIOHOMENS	TOTALRECEITA
2016	1210393.21	3232223.89	14614246.97
2017	2500857.97	7274421.87	17200423.35
2018	3800427.49	11135098.98	17593539.66
2019	473247.25	113834419.20	13535316.33
2020	473247.25	113834419.20	13022175.87