

# Otimização de Interrogações

Trabalho Prático nº 1



Mestrado Integrado em Engenharia Informática e  
Computação

Tecnologias de Bases de Dados

**Elementos do grupo:**

André Pires - 201207106 - ei12058@fe.up.pt  
João Bandeira - 201200615 - ei12022@fe.up.pt

Faculdade de Engenharia da Universidade do Porto  
Rua Roberto Frias, sn, 4200-465 Porto, Portugal

29 de Março de 2016

# Conteúdo

<b>1</b>	<b>Objetivo do trabalho</b>	<b>3</b>
<b>2</b>	<b>Criação de índices</b>	<b>3</b>
<b>3</b>	<b>Metodologia</b>	<b>3</b>
<b>4</b>	<b>Respostas às perguntas</b>	<b>4</b>
4.1	Pergunta 1 - Seleção . . . . .	4
4.1.1	Formulação SQL . . . . .	4
4.1.2	Tempos . . . . .	4
4.1.3	Plano de execução . . . . .	5
4.1.4	Resultados . . . . .	6
4.1.5	Conclusões . . . . .	6
4.2	Pergunta 2 - Agregação . . . . .	6
4.2.1	Formulação SQL . . . . .	6
4.2.2	Tempos . . . . .	7
4.2.3	Plano de execução . . . . .	8
4.2.4	Resultados . . . . .	11
4.2.5	Conclusões . . . . .	11
4.3	Pergunta 3 - Negação . . . . .	11
4.3.1	Formulação SQL . . . . .	11
4.3.2	Tempos . . . . .	12
4.3.3	Plano de execução . . . . .	12
4.3.4	Resultados . . . . .	14
4.3.5	Conclusões . . . . .	14
4.4	Pergunta 4 - Quantificação universal . . . . .	14
4.4.1	Formulação SQL . . . . .	14
4.4.2	Tempos . . . . .	16
4.4.3	Plano de execução . . . . .	17
4.4.4	Resultados . . . . .	19
4.4.5	Conclusões . . . . .	19

## 1 Objetivo do trabalho

Analisar os planos de execução de diferentes interrogações SQL a uma BD de teste. Avaliar o impacto da existência de índices e de estatísticas e do recurso a diferentes estratégias de estruturação das perguntas. O relatório a produzir deve, para cada situação, indicar a formulação da pergunta em SQL, a resposta obtida, as condições de execução do teste incluindo duração se relevante, o plano de execução e as conclusões a que chegou.[1]

## 2 Criação de índices

Para as tabelas "z", foram criados três índices cujo objetivo foi permitir uma maior eficiência de certas interrogações. Os índices criados foram:

- Btree na coluna "disciplina\_id" da tabela "zrespostas"

Foi criado um índice para a coluna relativa à disciplina, pois várias interrogações requerem agregação com a tabela das disciplinas, nesta mesma coluna.

- *Bitmap* na coluna "semestre\_id" da tabela "zrespostas"

Foi criado um índice para a coluna relativa ao semestre, pois várias interrogações esperam resultados de um semestre específico. Foi escolhido um índice *bitmap* devido às vantagens que o mesmo traz, e ao facto da cardinalidade da coluna semestre não ser muito elevada (existem respostas relativas a apenas alguns semestres).

- Btree nas colunas "semestre\_id" e "disciplina\_id" da tabela "zrespostas"

Este último índice foi criado para tentar melhorar a eficiência das interrogações que envolvem a agregação de resultados por disciplina, ao mesmo tempo que se filtra um semestre específico.

Para criação dos índices referidos foi utilizado o código SQL abaixo.

---

```
CREATE INDEX Z_RESPOSTAS_DISCIPLINA_IDX ON zrespostas (disciplina_id);

CREATE INDEX Z_RESPOSTAS_SEMESTRE_IDX ON zrespostas (semestre_id);

CREATE INDEX Z_RESP_DISC_SEMESTRE_IDX ON zrespostas (semestre_id,
disciplina_id);
```

---

## 3 Metodologia

Para cada uma das perguntas propostas no enunciado, o código SQL escrito foi corrido três vezes, sendo que o tempo de execução de cada uma delas foi registado. Sendo assim, os tempos apresentados abaixo correspondem a uma média dos tempos de cada conjunto de três execuções. Relativamente a estes tempos, apesar de serem apresentados para todas as perguntas, por vezes não serão muito significativos, uma vez que após a primeira execução de cada interrogação, frequentemente o servidor retornava resultados num tempo muito inferior, pois utiliza muito provavelmente informação em *cache*.

Foram ainda registados e apresentados abaixo os planos de execução e resultados de todas as interrogações.

## 4 Respostas às perguntas

### 4.1 Pergunta 1 - Seleção

#### 4.1.1 Formulação SQL

---

```
-- a
select count(*)
from xrespostas
where semestre_id = 21;

-- b
select count(*)
from xrespostas
Where disciplina_id = 1237;

-- c
select count(*)
from xrespostas, xsemestre
where xrespostas.semestre_id = xsemestre.semestre_id
and xsemestre.ano_lectivo = '2008/2009'
and xsemestre.semestre = '1S';
```

---

#### 4.1.2 Tempos

Alínea	X (ms)	Y (ms)	Z (ms)
a	130.67	130.33	29.67
b	132	132.67	40.67
c	309	285.33	42.67

Tabela 1: Resultados pergunta 1

### 4.1.3 Plano de execução

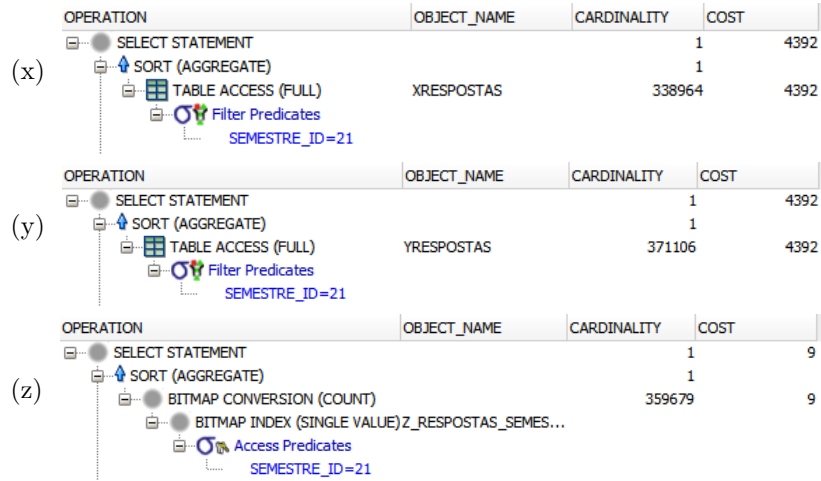


Figura 1: Plano de execução da alínea 1.a)

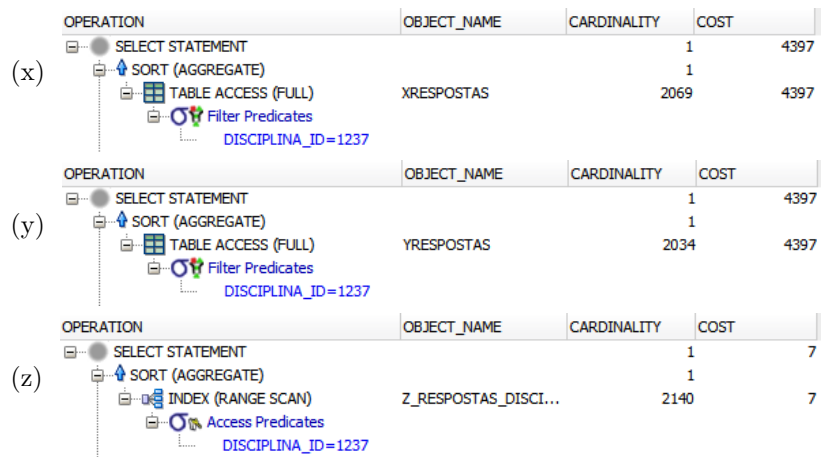


Figura 2: Plano de execução da alínea 1.b)

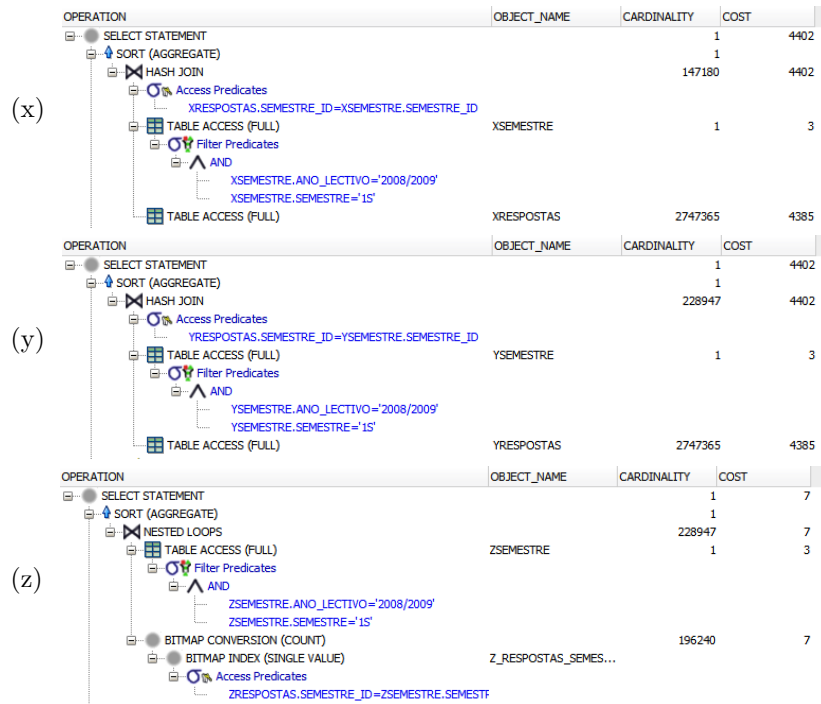


Figura 3: Plano de execução da alínea 1.c)

#### 4.1.4 Resultados

	a	b	c
Count(*)	351671	6525	351671

Tabela 2: Resultados pergunta 1

#### 4.1.5 Conclusões

Neste caso, nota-se claramente um resultado favorável, em tempo e em custo de execução, com a adição dos índices extra (tabelas z). No entanto, os índices "normais" (como *primary* e *foreign keys*) não têm qualquer influência na seleção.

## 4.2 Pergunta 2 - Agregação

### 4.2.1 Formulação SQL

```
-- a
select avg(xrespostas.resposta)
from xrespostas, xdisciplina
where xrespostas.disciplina_id = xdisciplina.disciplina_id
and xdisciplina.sigla = 'FPRO';

-- b
select * from (
```

```

select xdisciplina.disciplina_id, xdisciplina.sigla,
       avg(xrespostas.resposta) nota, count(*) nr
from xrespostas, xdisciplina
where xrespostas.semestre_id = 21
and xrespostas.disciplina_id = xdisciplina.disciplina_id
group by xdisciplina.disciplina_id, xdisciplina.sigla
having count(*) > 300
order by nota desc
)
where rownum <= 1;

-- c
CREATE VIEW x_media_mais_300_semestre_21 AS(
  SELECT disciplina_id, avg(resposta) nota_glob, COUNT(*) nr_respostas
  FROM xrespostas
  WHERE semestre_id = 21
  GROUP BY disciplina_id
  HAVING COUNT(*) > 300
);

CREATE VIEW x_med_pergunta_apreciacao_glob AS(
  SELECT xr.disciplina_id, avg(xr.resposta) nota_perg, COUNT(*)
    nr_respostas
  FROM xrespostas xr, xpergunta xp
  WHERE xr.semestre_id = 21
  AND xr.pergunta_id = xp.pergunta_id
  AND xp.nome LIKE 'Apreciacao Global'
  GROUP BY xr.disciplina_id
);

SELECT d.disciplina_id, d.sigla, med_glob.nota_glob, med_perg.nota_perg
FROM x_media_mais_300_semestre_21 med_glob,
     x_med_pergunta_apreciacao_glob med_perg, xdisciplina d
WHERE med_glob.disciplina_id = med_perg.disciplina_id
AND med_glob.disciplina_id = d.disciplina_id
AND nota_glob + 0.1 < nota_perg;

```

---

#### 4.2.2 Tempos

Alínea	X (ms)	Y (ms)	Z (ms)
a	299	306	18.33
b	296	262.67	256.67
c	520.67	368.33	318

Tabela 3: Resultados pergunta 2

### 4.2.3 Plano de execução

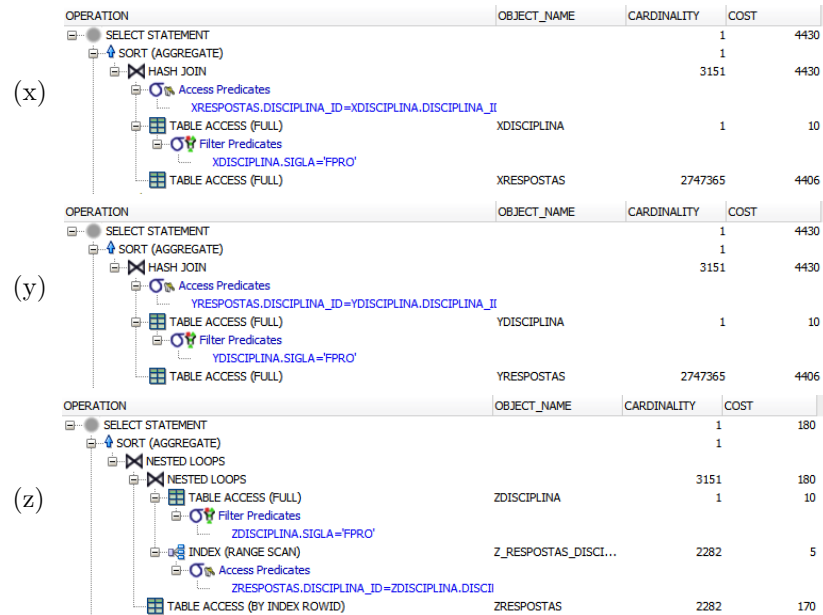


Figura 4: Plano de execução da alínea 2.a)



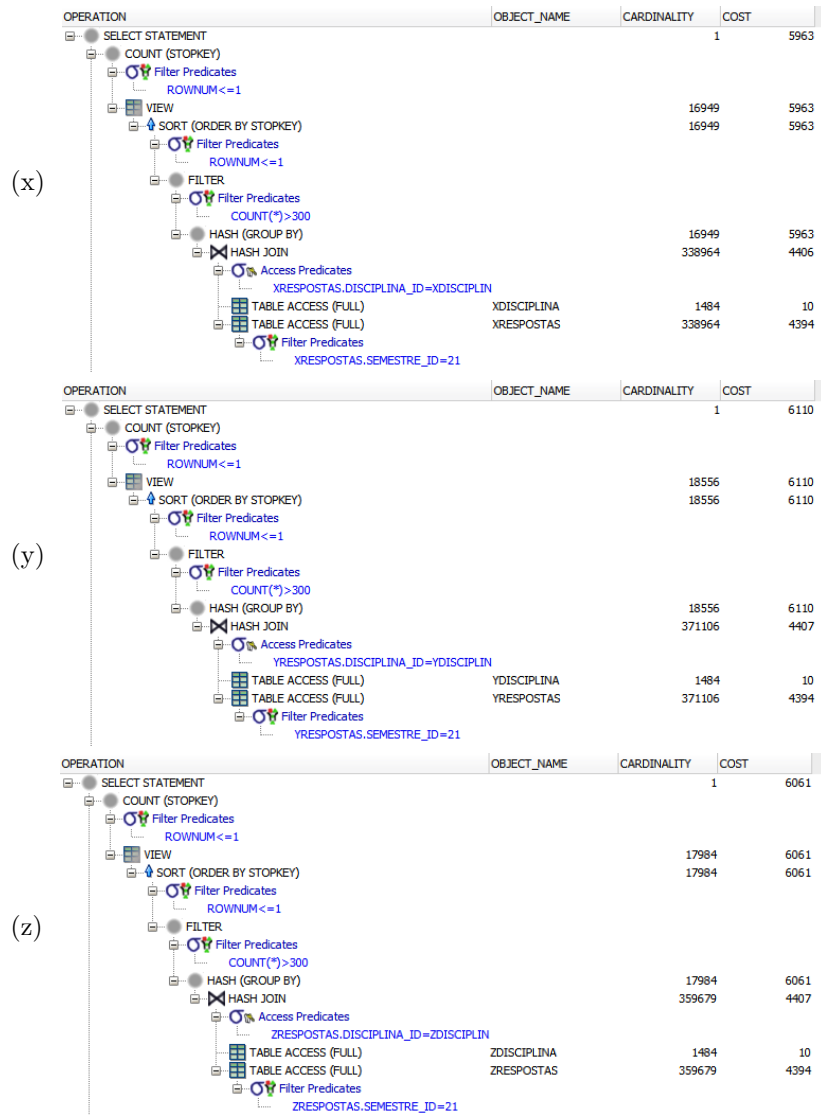


Figura 5: Plano de execução da alínea 2.b)

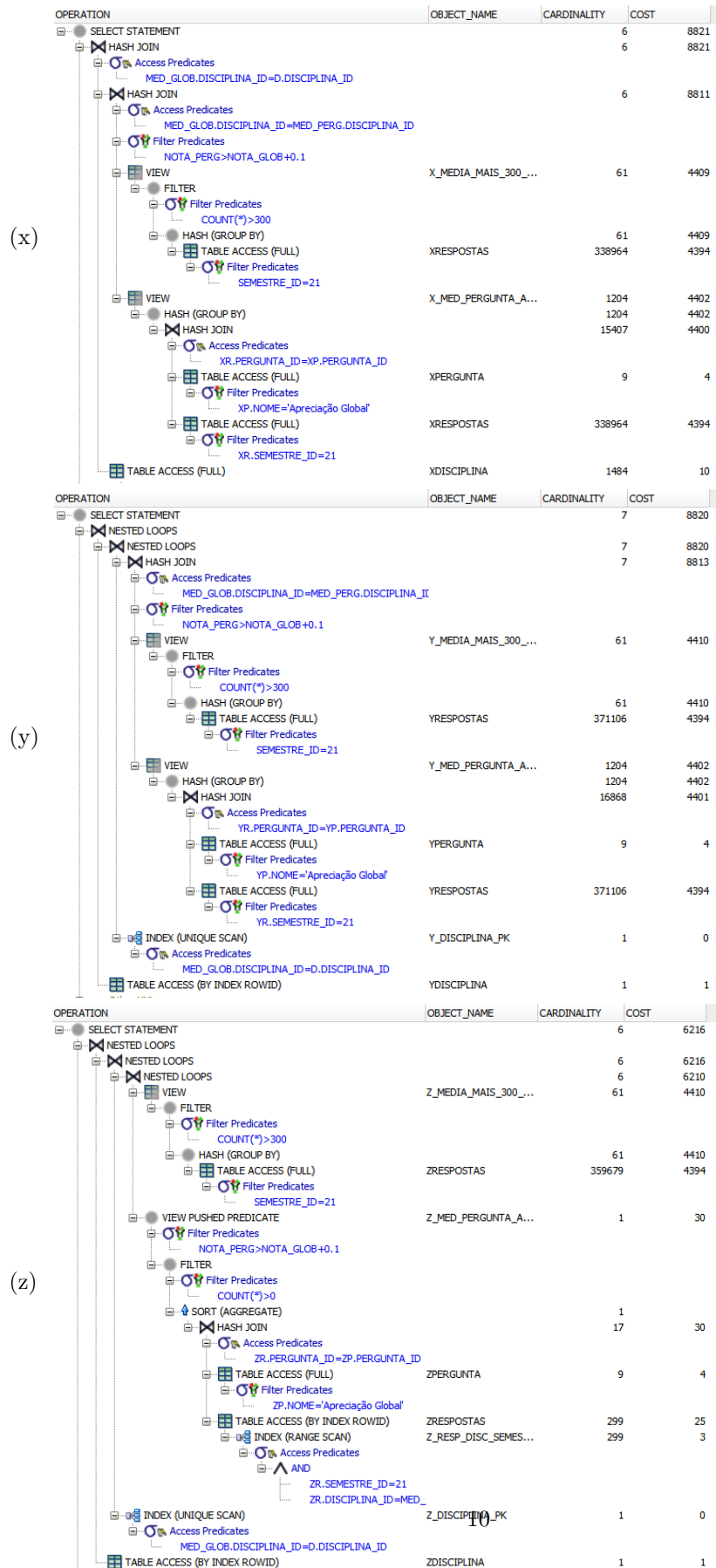


Figura 6: Plano de execução da alínea 2.c)

#### 4.2.4 Resultados

- Alínea a:  $\text{avg}() = 3,979617$
- Alínea b: disciplina\_id = 1949, sigla = SAD, nota = 4,519126
- Alínea c:

	DISCIPLINA_ID	SIGLA	NOTA_GLOB	NOTA_PERG
1	1172	MAT1	3,6309623430962343096234309623431	3,73684210526315789473684210526315789474
2	1661	TTC	4,056580565805658056580565805658057	4,166666666666666666666666666666667
3	1470	MICR	4,08639308855291576673866090712742980562	4,2
4	1511	QUIM1	3,97989949748743718592964824120603015075	4,1111111111111111111111111111111111
5	1510	FISI1	3,77577319587628865979381443298969072165	3,89473684210526315789473684210526315789
6	1907	AFHU	4,38709677419354838709677419354838709677	4,6363636363636363636363636363636364

Figura 7: Resultado 2c)

#### 4.2.5 Conclusões

Relativamente à alínea a), através da análise dos planos de execução pode verificar-se que para os conjuntos de tabelas x e y, as interrogações seguem uma execução semelhante. Apesar das tabelas y terem índices nas colunas que correspondem a chaves primárias, estas não são utilizadas na seleção. Já no terceiro plano de execução, relativo às tabelas z, observa-se que a interrogação tomou partido da existência de um índice na coluna 'disciplina\_id', apresentando um custo substancialmente inferior aos das tabelas x e y.

Já quanto à alínea c), é possível verificar diferenças entre todos os planos de execução. Na execução das tabelas y, pode verificar-se que esta tomou partido da existência de um índice na coluna 'disciplina\_id' da tabelas de disciplinas (chave primária). No que diz respeito às tabelas z, para além das melhorias que se observou em y em relação a x, a execução da interrogação foi ainda beneficiada pela existência do índice em múltiplas colunas (disciplina\_id e semestre\_id) na tabela de respostas. uma vez que esta interrogação envolve seleção de tuplos especificando o seu semestre, e agrupando-os por disciplina.

### 4.3 Pergunta 3 - Negação

#### 4.3.1 Formulação SQL

```
-- not in
select *
from xdisciplina
where disciplina_id not in(
    select disciplina_id
    from xrespostas
    group by disciplina_id
);

-- juncao externa e filtragem para nulo
select disc.sigla, disc.nome
from (
    select disciplina_id
    from xrespostas
```

```

group by disciplina_id
) resp FULL OUTER JOIN xdisciplina disc
ON resp.disciplina_id = disc.disciplina_id
where resp.disciplina_id is null;

```

---

#### 4.3.2 Tempos

Alínea	X (ms)	Y (ms)	Z (ms)
Baseada em <i>not in</i>	1020.67	788	448.67
Baseada em junção externa e filtragem para nulo	931.33	964.33	411

Tabela 4: Resultados pergunta 3, com otimizador = *CHOOSE*

Alínea	X (ms)	Y (ms)	Z (ms)
Baseada em <i>not in</i>	-	-	-
Baseada em junção externa e filtragem para nulo	419	419.67	418.33

Tabela 5: Resultados pergunta 3, com otimizador = *RULE*

#### 4.3.3 Plano de execução

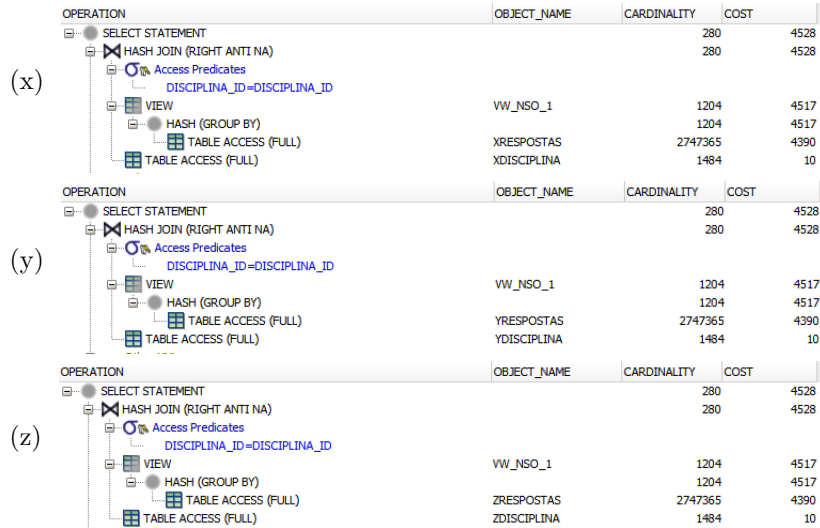


Figura 8: Plano de execução da alínea 3 *not in*, com otimizador *CHOOSE*

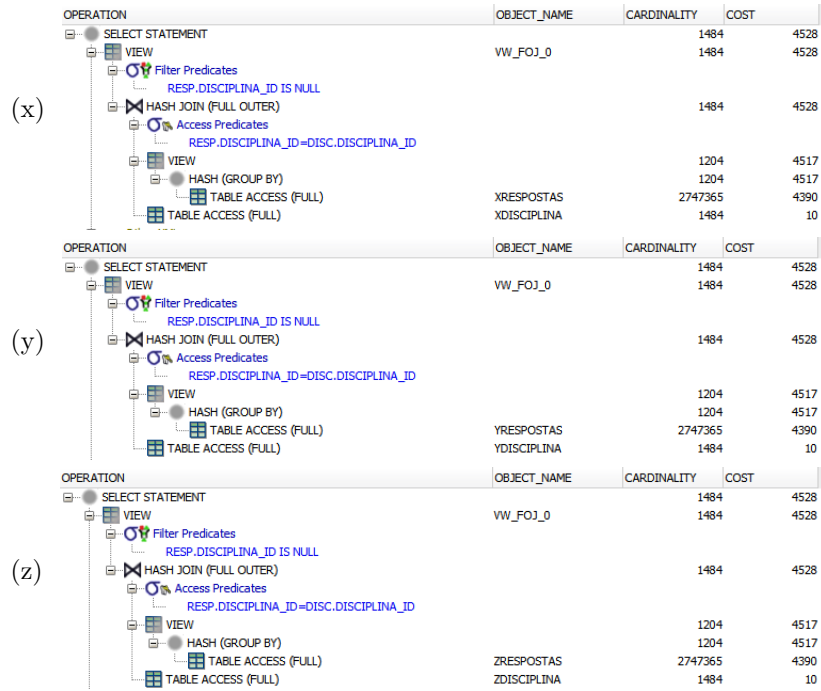


Figura 9: Plano de execução da alínea 3 *junção*, com otimizador *CHOOSE*

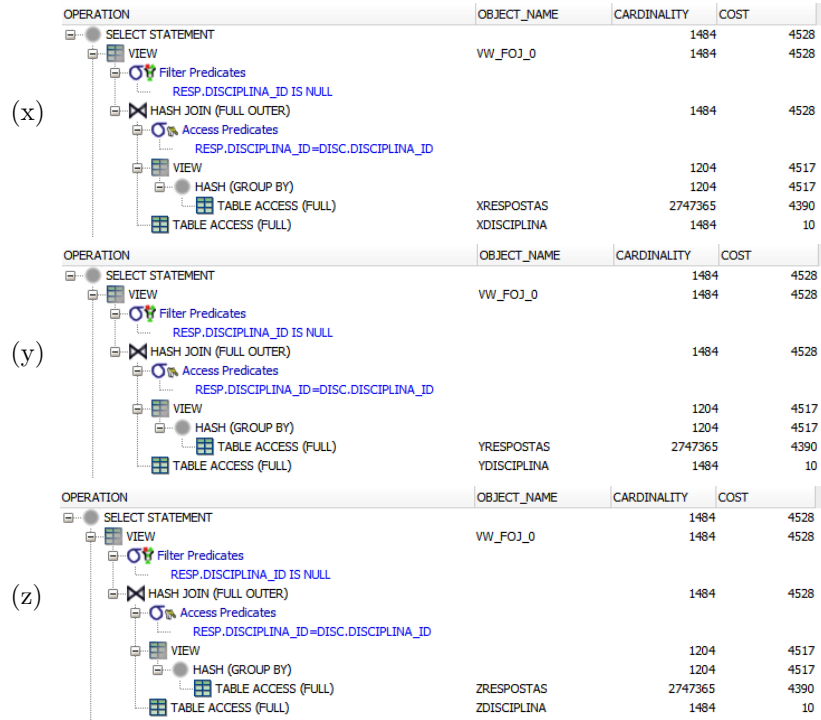


Figura 10: Plano de execução da alínea 3 *junção*, com otimizador *RULE*

#### 4.3.4 Resultados

	SIGLA	NOME
1	IE	Informação Empresarial
2	ASI	Análise de Sistemas de Informação
3	ETELE	Electrónica das Telecomunicações
4	PSD	Projecto de Sistemas Digitais
5	GA	Geologia Ambiental
6	PSA	Projecto de Sistemas de Automação
7	IMTM	Informática Médica e Telemedicina
8	LAUT	Laboratório de Automação
9	MIAGI	Metodologias de Investigação Aplicadas à Gestão de Informação
10	GI	Gestão de Informação
11	ARI	Armazenamento e Recuperação de Informação
12	SM	Sistemas Multimédia
13	CRC	Computadores e Redes de Comunicação
14	EMBS	Equipamento Médico e Biosensores
15	CR	Corrosão e Revestimentos
16	PSI	Produtos e Serviços de Informação
17	RM	Resistência dos Materiais
18	DRSIG	Deteção Remota e Sistemas de Informação Geográfica
19	BD	Bases de Dados
20	F I	Física I
21	SDO	Sistemas Dinâmicos e Optimização
22	CM	Caracterização de Materiais
23	ME	Métodos Estatísticos

Figura 11: Resultado 3

#### 4.3.5 Conclusões

Relativamente a esta pergunta, começamos por referir que não são apresentados resultados relativamente à execução da técnica "not in" com o otimizador "RULE" pois não se conseguiu que a interrogação executasse em tempo útil (após alguns minutos de espera, não tinham sido retornados resultados).

Analisando os planos de execução dos restantes casos, começando pelos casos onde foi utilizada a estratégia de junção externa e filtragem para nulo, verificamos que são semelhantes, não tendo a escolha do otimizador feito diferença alguma.

Já quanto à estratégia que utiliza "not in", e relativamente à estratégia anterior, apesar de requerer menos operações durante a execução, resulta num custo igual.

### 4.4 Pergunta 4 - Quantificação universal

#### 4.4.1 Formulação SQL

```
-- contagem
create view x_total_respostas as (
```

```

SELECT disciplina_id , count(*) total_respostas
FROM (
    SELECT disciplina_id, pergunta_id, count(*) c
    FROM xrespostas
    WHERE semestre_id = 21
    GROUP BY disciplina_id, pergunta_id
)
GROUP BY disciplina_id
);

create view x_respostas_5 as (
    SELECT disciplina_id , count(*) respostas_5
    FROM (
        SELECT disciplina_id, pergunta_id, count(*)
        FROM xrespostas
        WHERE semestre_id = 21
        AND resposta = 5
        GROUP BY DISCIPLINA_ID, pergunta_id
    )
    GROUP BY disciplina_id
);

SELECT d.disciplina_id, d.sigla, t.total_respostas, r.respostas_5
FROM x_total_respostas t, x_respostas_5 r, xdisciplina d
WHERE t.disciplina_id = r.disciplina_id
AND t.disciplina_id = d.disciplina_id
AND t.total_respostas = r.respostas_5;

-- dupla negacao
create view x_disciplinas_inq_21 as (
    SELECT disciplina_id
    FROM xrespostas
    WHERE semestre_id = 21
    GROUP BY disciplina_id
);

SELECT db.disciplina_id, db.sigla
FROM x_disciplinas_inq_21 da, xdisciplina db
WHERE da.disciplina_id = db.disciplina_id
AND da.disciplina_id NOT IN (
    SELECT disciplina_id
    FROM xrespostas ra
    WHERE semestre_id = 21
    AND NOT EXISTS (
        SELECT rb.disciplina_id, rb.pergunta_id
        FROM xrespostas rb
        WHERE rb.disciplina_id = ra.disciplina_id
        AND rb.pergunta_id = ra.pergunta_id
        AND rb.semestre_id = 21
        AND rb.resposta = 5
        GROUP BY disciplina_id, pergunta_id
    )
)
GROUP BY disciplina_id, pergunta_id
);

```

---

#### 4.4.2 Tempos

<b>Alínea</b>	<b>X (ms)</b>	<b>Y (ms)</b>	<b>Z (ms)</b>
Estratégia de dupla negação	336.67	365.67	294
Estratégia de contagem	496	496	393.33

Tabela 6: Resultados pergunta 4



#### 4.4.3 Plano de execução

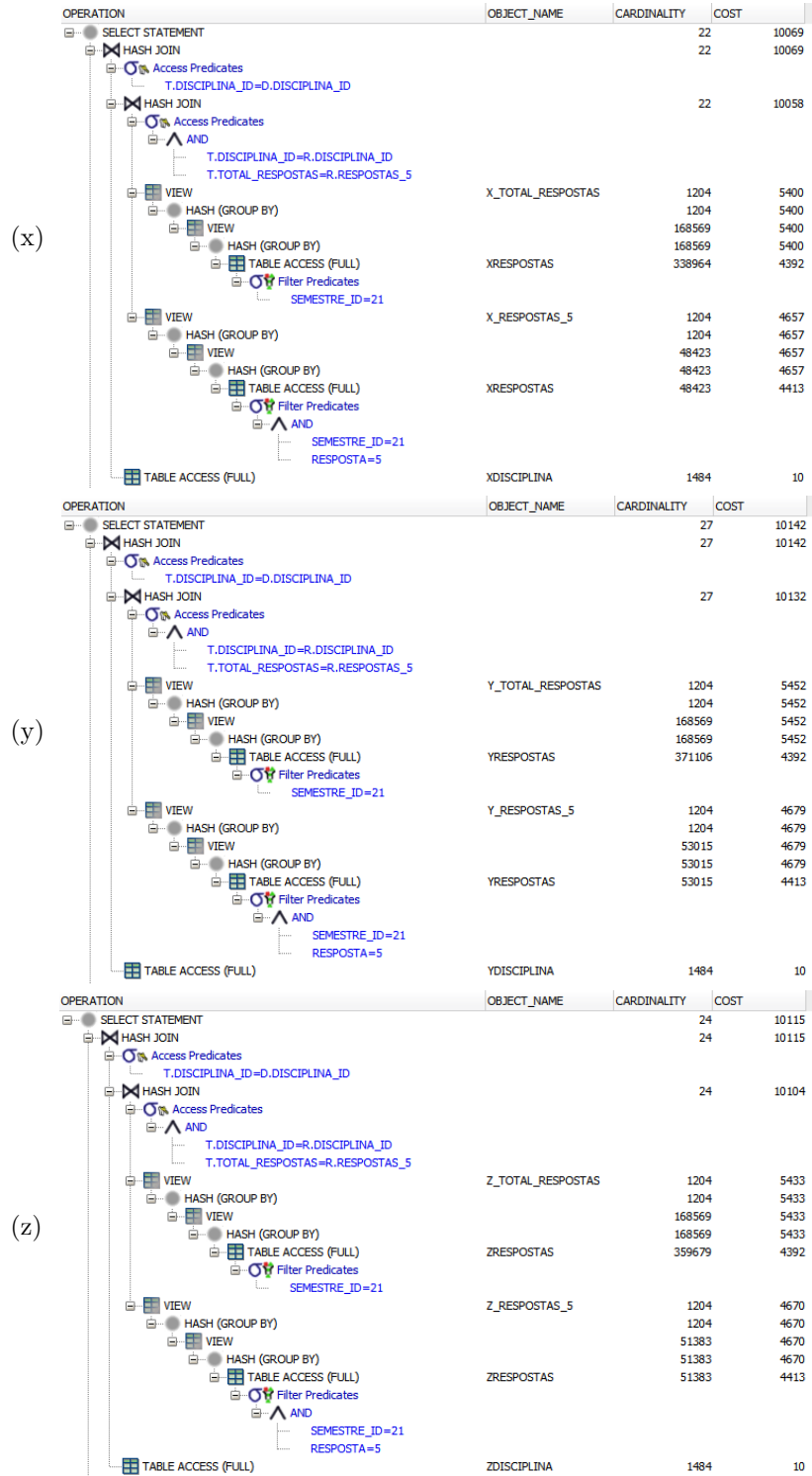


Figura 12: Plano de execução da alínea 4 *contagem*

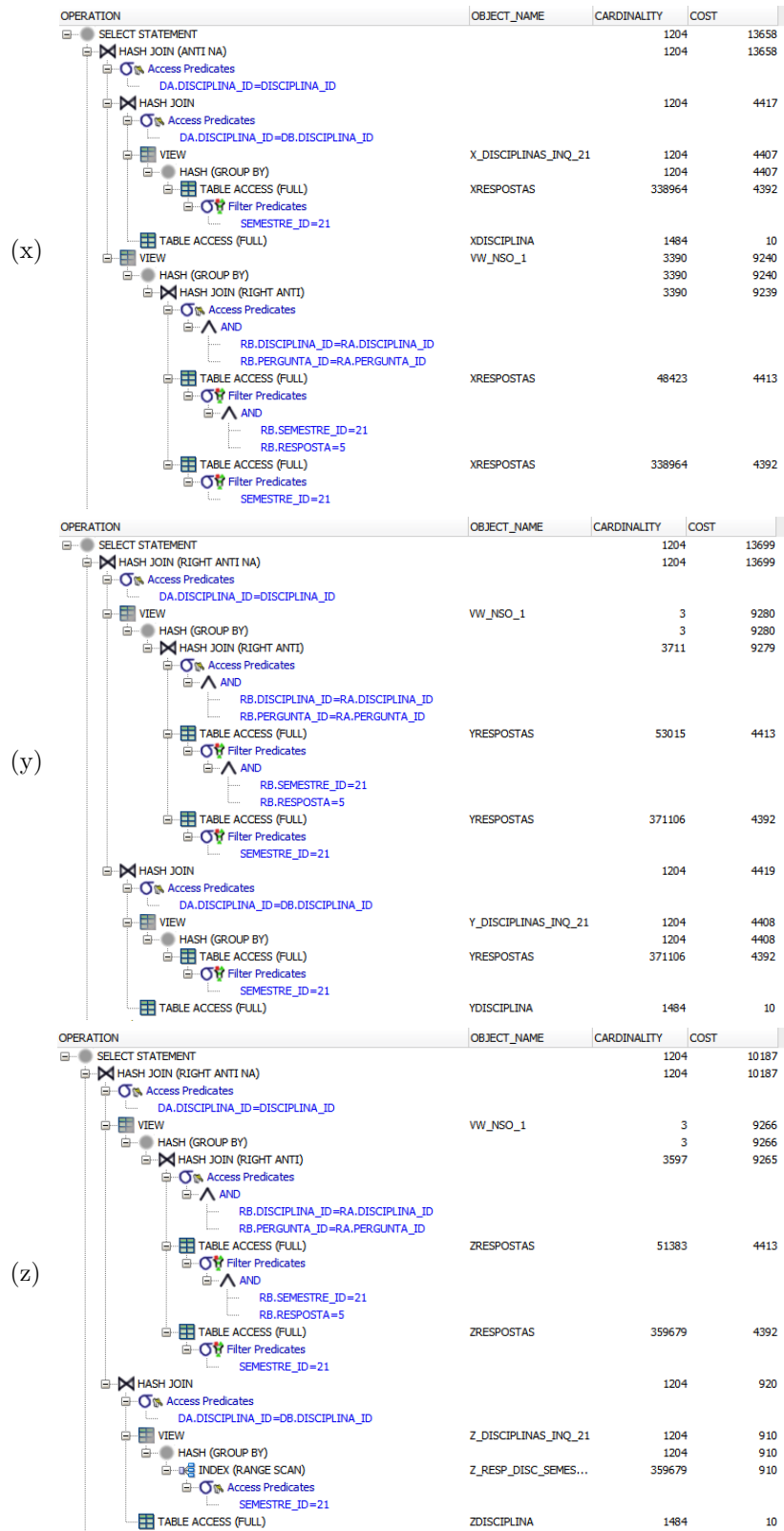


Figura 13: Plano de execução da alínea 4 *dupla negação*

#### 4.4.4 Resultados

	DISCIPLINA_ID	SIGLA
1	1235	ALGE
2	1236	AMAT
3	1237	FPRO
4	1238	IOCO
5	1239	IBDA
6	1245	AEDA
7	1265	AIAD
8	1266	GEMP
9	1267	AMAT1
10	1268	ALGE
11	1269	PROG1
12	1270	LSDI
13	1387	LESO
14	1388	SDIS
15	1393	MCM II
16	1402	EI
17	1403	TC
18	1405	IO
19	1124	GOPE
20	1144	RTSEE
21	1156	TPRE
22	1158	EESO
23	1688	ALGA

Figura 14: Resultado 4

#### 4.4.5 Conclusões

Na execução das interrogações escritas para esta pergunta, verificou-se que existem diferenças na execução de cada uma das estratégias (contagem e dupla negação), na medida em que, para os conjuntos de tabelas x e y, existe uma diferença de custo relativamente significativa entre ambas, com vantagem para a estratégia que envolve contagem.

No que diz respeito ao conjunto de tabelas z, a diferença entre as duas estratégias não foi significativa, e até desprezável, uma vez que na estratégia de dupla negação o sistema conseguiu tomar partido da existência de um índice nas colunas disciplina.id e semestre.id, reduzindo o custo total da interrogação.

## Referências

- [1] Gabriel David. Otimização de interrogações. Enunciado do projeto, [https://moodle.up.pt/pluginfile.php/102650/mod\\_resource/content/0/37-TrabOpt3.pdf](https://moodle.up.pt/pluginfile.php/102650/mod_resource/content/0/37-TrabOpt3.pdf), 2016.