# INSTITUTO INFNET PÓS EM ENGENHARIA DE DADOS

Juliana Opis Mikosz

Hadoop

## Introdução:

Este trabalho tem como objetivo realizar uma análise de dados em dois conjuntos de dados distintos: Filmes da Netflix e Transações de Compras.

Para a execução desse trabalho foi necessária a criação de um cluster no dataproc, em seguida a importação dos datasets escolhidos para um bucket dentro do GCP. Após, foi feita a cópia desses arquivos para o HDFS, criação de databases e tabelas no Hive através do Beeline. Os datasets armazenados no HDFS foram então injetados nas tabelas correspondentes no Hive. Por último algumas consultas para a análise de dados e obtenção de insights relevantes.

## Conjuntos de dados:

#### Conjunto de Dados - Filmes da Netflix:

O primeiro conjunto de dados foi obtido a partir de fontes públicas no Kaggle. Ele inclui informações sobre filmes disponíveis na plataforma Netflix, como o ano de lançamento e as avaliações dos filmes. Estes dados serão utilizados para analisar e explorar a evolução das avaliações dos filmes ao longo do tempo, bem como para identificar padrões de qualidade com base nas avaliações dos usuários da plataforma.

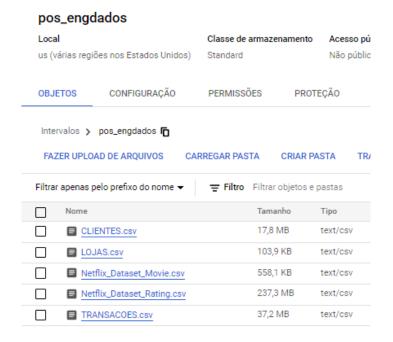
Fonte: https://www.kaggle.com/datasets/rishitjavia/netflix-movie-rating-dataset

#### Conjunto de Dados - Transações de Compras:

O segundo conjunto de dados foi extraído do banco de dados da nossa empresa, mas todas as informações sensíveis foram anonimizadas e substituídas por dados fictícios para garantir a privacidade e a segurança dos dados. Este conjunto de dados contém informações sobre transações de compras realizadas pela empresa, incluindo detalhes como identificação da loja, cliente, valor da compra e datas das transações. Será realizada uma análise detalhada dessas transações para identificar tendências de compra, comportamentos dos clientes e padrões sazonais.

## Passos do projeto:

1. Exportação dos arquivos do Bucket para o HDFS:



#### 2. Criação da pasta datasets:

Comando: mkdir datasets

#### 3. Importação dos arquivos do bucket:

Comando: gcsfuse pos\_engdados datasets

Datasets criados na pasta /datasets:

```
juliana_mikosz@cluster-1aa7-m:~/datasets$ 1s -1a

total 299981
-rw-r--r- 1 juliana_mikosz juliana_mikosz 18657607 Aug 31 21:45 CLIENTES.csv
-rw-r--r- 1 juliana_mikosz juliana_mikosz 106443 Aug 31 21:45 LOJAS.csv
-rw-r--r- 1 juliana_mikosz juliana_mikosz 571494 Aug 31 21:45 Netflix_Dataset_Movie.csv
-rw-r--r- 1 juliana_mikosz juliana_mikosz 248836313 Aug 31 21:46 Netflix_Dataset_Rating.csv
-rw-r--r- 1 juliana_mikosz juliana_mikosz 39007358 Aug 31 21:45 TRANSACOES.csv
```

#### 4. Copia os arquivos locais para o HDFS:

Comando: hdfs dfs -put \*.csv /user/juliana\_mikosz/datasets

#### 5. Conferência se os arquivos foram copiados:

Comando: hdfs dfs -ls /user/juliana\_mikosz/datasets

#### Resultado:

```
juliana_mikosz@cluster-1aa7-m:~$ hdfs dfs -ls /user/juliana_mikosz/datasets
Found 5 items
          1 juliana_mikosz hadoop 18657607 2023-08-31 21:56 /user/juliana_mikosz/datasets
-rw-r--r--
/CLIENTES.csv
-rw-r--r--
           1 juliana_mikosz hadoop
                                      106443 2023-08-31 21:56 /user/juliana_mikosz/datasets
/LOJAS.csv
          1 juliana_mikosz hadoop 571494 2023-08-31 21:56 /user/juliana_mikosz/datasets
-rw-r--r--
/Netflix Dataset Movie.csv
-rw-r--r- 1 juliana_mikosz hadoop 248836313 2023-08-31 21:56 /user/juliana_mikosz/datasets
/Netflix_Dataset_Rating.csv
-rw-r--r- 1 juliana_mikosz hadoop 39007358 2023-08-31 21:56 /user/juliana_mikosz/datasets
/TRANSACOES.csv
```

#### Hive:

#### 1. Conexão beeline:

Comando: beeline -u jdbc:hive2://localhost:10000/default -n juliana\_mikosz@cluster-1aa7-m -d org.apache.hive.jdbc.HiveDriver

#### Conexão realizada:

```
juliana_mikosz@cluster-1aa7-m:~$ beeline -u jdbc:hive2://localhost:10000/default -n juliana_mi
kosz@cluster-1aa7-m -d org.apache.hive.jdbc.HiveDriver
Connecting to jdbc:hive2://localhost:10000/default
Connected to: Apache Hive (version 3.1.3)
Driver: Hive JDBC (version 3.1.3)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 3.1.3 by Apache Hive
0: jdbc:hive2://localhost:10000/default>
```

#### 2. Criação dos databases e tabelas

Comando: create database if not exists vendas comment "Database com informações de clientes, vendas e lojas";

Comando: create database if not exists netflix comment "Database com informações de filmes da Netflix";

Databases criados:

```
0: jdbc:hive2://localhost:10000/default> create database if not exists vendas comment "Databas
e com informações de clientes, vendas e lojas";
INFO : Compiling command(queryId=hive_20230831234118_255f6fb8-e8af-4e88-a7b6-2be1e2687d4e): c
reate database if not exists vendas comment "Database com informações de clientes, vendas e lo
jas"
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldSchemas:null, properties:null)
      : Completed compiling command(queryId=hive_20230831234118_255f6fb8-e8af-4e88-a7b6-2be1e2
687d4e); Time taken: 0.019 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
     : Executing command(queryId=hive_20230831234118_255f6fb8-e8af-4e88-a7b6-2bele2687d4e): c
reate database if not exists vendas comment "Database com informações de clientes, vendas e lo
jas"
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId=hive_20230831234118_255f6fb8-e8af-4e88-a7b6-2be1e2
687d4e); Time taken: 0.033 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
No rows affected (0.068 seconds)
0: jdbc:hive2://localhost:10000/default> create database if not exists netflix comment "Databa
se com informações de filmes da Netflix";
INFO : Compiling command(queryId=hive 20230831234139 c6034228-8c9f-475c-9020-7ba547d159a4): c
reate database if not exists netflix comment "Database com informações de filmes da Netflix"
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldSchemas:null, properties:null)
     : Completed compiling command(queryId=hive 20230831234139 c6034228-8c9f-475c-9020-7ba547
d159a4); Time taken: 0.024 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20230831234139_c6034228-8c9f-475c-9020-7ba547d159a4): c
reate database if not exists netflix comment "Database com informações de filmes da Netflix"
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId=hive_20230831234139_c6034228-8c9f-475c-9020-7ba547
d159a4); Time taken: 0.043 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
No rows affected (0.083 seconds)
```

```
0: jdbc:hive2://localhost:10000/default> show databases;
INFO : Compiling command(queryId=hive 20230831234220 a422e91d-1928-4238-a930-8a607412e0cf): s
how databases
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:database name, type:strin
g, comment:from deserializer)], properties:null)
INFO : Completed compiling command(queryId=hive_20230831234220_a422e91d-1928-4238-a930-8a6074
12e0cf); Time taken: 0.018 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive 20230831234220 a422e91d-1928-4238-a930-8a607412e0cf): s
how databases
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId=hive 20230831234220 a422e91d-1928-4238-a930-8a6074
12e0cf); Time taken: 0.005 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
| database_name |
| default
| netflix
| testejuliana
| vendas
```

#### Dicionário de dados tabelas | Database Vendas:

#### 1. Tabela clientes

id (INT): Um identificador único para o cliente.

nome (STRING): O nome do cliente.

ativo (INT): Indica se o cliente está ativo ou não.

genero (STRING): O gênero do cliente.

city (STRING): A cidade em que o cliente reside.

#### 2. Tabela transacoes

Id (INT): Um identificador único para a transação.

loja (INT): O identificador único da loja onde a transação ocorreu.

cliente (INT): O identificador único do cliente que fez a transação.

valor (INT): O valor da transação.

data (DATE): A data e hora em que a transação ocorreu.

#### 3. Tabela lojas

Id (INT): Um identificador único para a loja.

ativo (INT): Indica se a loja está ativa ou não.

cidade (STRING): A cidade onde a loja está localizada.

data instalação (DATE): A data de instalação da loja.

#### Criação das tabelas | Database vendas:

#### 1. Criação da tabela clientes:

CREATE TABLE clientes (id INT, nome STRING, ativo INT, genero STRING, city STRING)

**ROW FORMAT DELIMITED** 

FIELDS TERMINATED BY ';'

STORED AS TEXTFILE:

#### 2. Ingestão dos dados na tabela criada:

LOAD DATA INPATH '/user/juliana\_mikosz/datasets/CLIENTES.csv' overwrite into table clientes;

#### 3. Criação da tabela transacoes:

CREATE TABLE transacoes (Id INT, loja INT, cliente INT, valor INT, datahora DATE)

**ROW FORMAT DELIMITED** 

FIELDS TERMINATED BY ':'

STORED AS TEXTFILE;

#### 4. Ingestão dos dados na tabela criada:

LOAD DATA INPATH '/user/juliana\_mikosz/datasets/TRANSACOES.csv' overwrite into table transacoes;

#### 5. Criação da tabela lojas:

CREATE TABLE lojas (Id INT, ativo INT, cidade STRING, data\_instalacao DATE) ROW FORMAT DELIMITED FIELDS TERMINATED BY ';' STORED AS TEXTFILE;

#### 6. Ingestão dos dados na tabela criada:

LOAD DATA INPATH '/user/juliana\_mikosz/datasets/LOJAS.csv' overwrite into table lojas;

#### Dicionário de dados tabelas | Database Vendas:

#### 1. Tabela filme

Id (INT): Um identificador único para o filme. year (INT): O ano de lançamento do filme. name (STRING): O nome do filme.

#### 2. Tabela avaliação

Id (INT): Um identificador único para a avaliação.
rating (INT): A classificação atribuída ao filme na avaliação.
movie id (INT): O identificador único do filme ao qual esta avaliação está associada.

#### Criação das tabelas | Database netflix:

#### 1. Criação da tabela filme:

CREATE TABLE filme (movie\_id INT, year INT, name STRING)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE;

#### 2. Ingestão dos dados na tabela criada:

LOAD DATA INPATH '/user/juliana\_mikosz/datasets/Netflix\_Dataset\_Movie.csv' overwrite into table filme:

#### 3. Criação da tabela avaliacao:

CREATE TABLE avaliacao (user\_id INT, rating INT, movie\_id INT) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' STORED AS TEXTFILE;

#### 4. Ingestão dos dados na tabela criada:

LOAD DATA INPATH '/user/juliana\_mikosz/datasets//Netflix\_Dataset\_Rating.csv' overwrite into table avaliacao;

## **Consultas Hive:**

## 1. Database Vendas:

## a) Qual o total de transações por cliente?

Consulta realizada:
SELECT c.id, c.nome, COUNT(\*) AS total\_transacoes
FROM clientes c
LEFT JOIN transacoes t ON c.id = t.cliente
GROUP BY c.id, c.nome
ORDER BY total\_transacoes DESC
LIMIT 20;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1 Map 5 Reducer 2 Reducer 3 Reducer 4	container container container	SUCCEEDED SUCCEEDED SUCCEEDED SUCCEEDED SUCCEEDED	1 1 3 1	1 1 3 1 1	0 0 0 0 0	0 0 0 0 0	0 0 0 0 0	0 0 0 0 0
VERTICES: 05/05	[	 	>>]	100% ELAPS	ED TIME:	23.04 s	<del></del>	

+-		-+		+-		-+
1	c.id	-1	c.nome	- 1	total_transacoes	- 1
+		-+		+-		-+
1	450405	1	Gabriel	-1	2182	-1
1	505132	-1	Manuela	- 1	1097	- 1
Τ	444612	-1	Manuela	- 1	1007	- 1
1	405123	-1	Afonso	-1	786	- 1
1	432171	-1	Pedro	-1	747	- 1
1	450426	-1	Alice	-1	716	-1
1	419959	-1	Lucas	-1	664	- 1
1	525734	-1	Guilherme	-1	605	- 1
Ι	559536	-1	Laura	-1	585	- 1
Τ	601902	-1	Marina	-1	538	- 1
1	438969	-1	Gabriel	-1	456	-1
ī	498439	-1	Caio	-1	435	-1
Ι	512916	-1	Muriel	-1	340	-1
ī	280302	-1	Bianca	-1	316	-1
Ι	455609	-1	Isaac	-1	253	- 1
ı	513277	1	Miguel	ı	247	1
ı	284187	1	Gabriel	-	224	1
I	517632	ı	Eduarda	-1	212	
١	405468	1	Alice	ı	206	I
١	462861	1	Felipe	-1	204	I
+-		-+		-+-		-+

## b) Qual o total de vendas por cidade?

Consulta realizada:
SELECT I.cidade, COUNT(t.id) AS total\_vendas
FROM lojas I
LEFT JOIN transacoes t ON I.id = t.loja
GROUP BY I.cidade
ORDER BY total\_vendas DESC;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	0	0	0	0
Map 5	container	SUCCEEDED	1	1	0	0	0	0
Reducer 2	container	SUCCEEDED	2	2	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0
Reducer 4	container	SUCCEEDED	1	1	0	0	0	0
VERTICES: 05/05			>>1	 100% ELAPS	ED TIME:	14.08 s		

_	
l l.cidade	+
+	+
São Paulo	
Curitiba	113490
Rio de Janeiro	1 35798
Belo Horizonte	I 26875
Guarulhos	I 26062 I
Blumenau	i 23599 i
Brasília	22371
Salvador	18830
Porto Alegre	18027
Londrina	16353
Maceió	13871
Jundiaí	12941
Vitória	12396
Osasco	11622
Santo André	10918
São José dos Campos	10560
Cuiabá	10104
Campo Grande	10019
Santos	9823
Teresina	8531
Vila Velha	8412
Uberlândia	8358
Barueri	8282
Florianópolis	8247
São José	8160
No <b>v</b> o Hamburgo	8030
Cascavel	7588
Nova Lima	6783
João Pessoa	6462
Bauru	5829
Joinville	5700
Araçatuba	5667
Diadema	5345
Niterói	5331
Campinas	5276
Maringá	5263
Recife	5073

## c) Qual o valor total de compras por cliente?

Consulta realizada:
SELECT c.id, c.nome, SUM(t.valor) AS valor\_total\_compras
FROM clientes c
LEFT JOIN transacoes t ON c.id = t.cliente
GROUP BY c.id, c.nome
ORDER BY valor\_total\_compras DESC
LIMIT 20;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	0	 0	 0	0
Map 5	container	SUCCEEDED	1	1	0	0	0	0
Reducer 2	container	SUCCEEDED	3	3	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0
Reducer 4	container	SUCCEEDED	1	1	0	0	0	0
VERTICES: 05/05	[			100% ELAPS	ED TIME:			·

	oosia.		
+			-+
1	c.id	c.nome   valor_total_compras	1
+			-+
1	450405	Gabriel   194545	1
1	505132	Manuela   99883	1
1	444612	Manuela   88969	1
1	405123	Afonso   68857	1
1	432171	Pedro   63592	1
1	450426	Alice   61619	1
1	419959	Lucas   59134	1
1	525734	Guilherme   55120	1
1	559536	Laura   53434	1
1	601902	Marina   49801	1
1	438969	Gabriel   40300	1
1	498439	Caio   37255	1
1	512916	Muriel   30425	1
1	280302	Bianca   27993	1
1	455609	Isaac   22594	1
1	513277	Miguel   20546	I
1	284187	Gabriel   19302	I
1	405468	Alice   18211	I
1	462861	Felipe   17911	I
1	517632	Eduarda   17306	1
+			-+
•		· · · · · · · · · · · · · · · · · · ·	•

## d) Qual a quantidade de clientes ativos por loja?

Consulta realizada:
SELECT I.id, I.cidade, COUNT(DISTINCT c.id) AS clientes\_ativos
FROM lojas I
LEFT JOIN transacoes t ON I.id = t.loja
LEFT JOIN clientes c ON t.cliente = c.id
WHERE c.ativo = 1
GROUP BY I.id, I.cidade
ORDER BY clientes\_ativos DESC;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	 0	 0	 0	 0
Map 2		SUCCEEDED	1	1	0	0	0	0
Map 3	container	SUCCEEDED	1	1	0	0	0	0
Reducer 4	container	SUCCEEDED	1	1	0	0	0	0
Reducer 5	container	SUCCEEDED	1	1	0	0	0	0
				100% ELAPS	ED TIME:	21.75 s		

1		
1.id	1.cidade	clientes ativos
+		-+
1 3225	São Paulo	1 1263
1 1555		I 819 I
1 2245	Recife	801
1821	São José	1 764
1001	Curitiba	740
1741	Santos	i 732 i
1707	Salwador	696
1823	Londrina	660
2432	São Paulo	641
2655	São Paulo	638
1043	Curitiba	599
1411	Jundiaí	589
1112	Campo Grande	589
3252	Curitiba	588
2565	São Paulo	563
2868	São Paulo	547
3212	Brasília	540
1039	Curitiba	519
1014	Curitiba	516
1231	Belo Horizonte	514
1862	Londrina	512
3337	Curitiba	495
3249	Juiz de Fora	495
2985	Brasília	491
1893	Votorantim	483
1746		463
	Curitiba	456
	São Paulo	454
1121	<u>-</u> -	438
2968		436
2407		422
1052		417
2175		411
3105		409
1095		409
2073		408
2767		408
2364		407
1315		403
2132	Porto Alegre	400

## e) Qual a quantidade de clientes ativos por cidade?

Consulta realizada:
SELECT I.cidade, COUNT(DISTINCT c.id) AS clientes\_ativos
FROM lojas I
LEFT JOIN transacoes t ON I.id = t.loja
LEFT JOIN clientes c ON t.cliente = c.id
WHERE c.ativo = 1
GROUP BY I.cidade
ORDER BY clientes\_ativos DESC;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	 0	 0	0	0
Map 2	container	SUCCEEDED	1	1	0	0	0	0
Map 3	container	SUCCEEDED	1	1	0	0	0	0
Reducer 4	container	SUCCEEDED	1	1	0	0	0	0
Reducer 5	container	SUCCEEDED	1	1	0	0	0	0
VERTICES: 05/05	[		>>]	100% ELAPS	ED TIME:	21.69 s		

+	++
l.cidade	clientes_ativos
+	++
São Paulo	61377
Curitiba	22951
Rio de Janeiro	6145
Guarulhos	5520
Belo Horizonte	5433
Brasília	5390
Sal <b>v</b> ador	4718
Blumenau	3380
Londrina	3248
Santos	2846
Cuiabá	2796
Porto Alegre	2694
Maceió	2686
Jundiaí	2594
Santo André	2518
Vitória	2335
Osasco	2316
Campo Grande	2273
Teresina	2142
São José dos Campos	1898
Cascavel	1841
Barueri	1720
Uberlândia	1677
Nova Lima	1508
Niterói	1395
São José	1334
Novo Hamburgo	1304
Mogi das Cruzes	1282
Vila Velha	1281
Florianópolis	1240
Recife	1224
Diadema	1218
Bauru	1186
Araçatuba	1179
Campinas	1176
João Pessoa	1166
Biguaçu	1139
Caruaru	1110
São Bernardo do Campo	1070

## f) Quais são as 3 lojas com o maior valor de venda em cada mês de 2023?

```
Consulta realizada:
WITH VendasPorLojaPorMes AS (
  SELECT
    I.id AS id_loja,
    I.cidade AS cidade,
    EXTRACT(MONTH FROM t.datahora) AS mes,
    SUM(t.valor) AS total_vendas,
    DENSE_RANK() OVER (PARTITION BY EXTRACT(MONTH FROM t.datahora)
ORDER BY SUM(t.valor) DESC) AS ranking
  FROM lojas I
  LEFT JOIN transacoes t ON I.id = t.loja
  WHERE I.id IS NOT NULL AND EXTRACT(MONTH FROM t.datahora) IS NOT
NULL
  GROUP BY EXTRACT(MONTH FROM t.datahora), I.id, I.cidade
SELECT id_loja, cidade, mes, total_vendas
FROM VendasPorLojaPorMes
```

# WHERE ranking <= 3 AND mes IS NOT NULL ORDER BY mes, ranking;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 5 Map 1 Reducer 2 Reducer 3 Reducer 4	container container container	SUCCEEDED SUCCEEDED SUCCEEDED SUCCEEDED SUCCEEDED	1 1 2 2 1	1 1 2 2 1	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0
VERTICES: 05/05	[	 	>>]	100% ELAPS	ED TIME:	14.95 s		

_								
+-		-+-		-+-		+-		-+
ı	id_loja	ı	cidade	ı	mes	ı	total_vendas	- 1
+-		-+-		-+-		+-		-+
L	1555	Ι	Guarulhos	П	1	Т	240193	-1
L	2432	Ι	São Paulo		1	Т	232895	-1
L	3248	1	São Paulo	-1	1	Т	224702	-1
L	3248	1	São Paulo		2	1	175888	-1
L	1821	1	São José	-1	2	1	160414	-1
L	2767	1	São Paulo	-1	2	Τ	138667	-1
L	1043	1	Curitiba	-1	3	1	7160	-1
L	2320	1	Suzano	-1	3	1	7098	-1
L	2655	1	São Paulo	-1	3	Τ	6559	-1
L	2655	1	São Paulo	-1	4	Τ	4888	-1
L	2013	1	Brasília	-1	4	Τ	3592	1
L	2432	1	São Paulo	-1	4	Τ	3547	1
Ī	2407	Τ	Salvador	-1	5	Τ	2663	1
Ī	1707	1	Salvador	-1	5	Τ	2593	1
Ī	1043	1	Curitiba	-1	5	Τ	2561	1
ī	1741	1	Santos	-1	6	Τ	2389	1
ī	2013	Τ	Brasília	-1	6	Τ	2057	1
ī	2407	1	Salvador	-1	6	Τ	2029	1
ī	1741	1	Santos	-1	7	Τ	3113	1
Ī	2655	1	São Paulo	-1	7	Τ	2238	1
I	2013	I	Brasília	I	7	1	2213	1
I	2132	I	Porto Alegre	I	8	1	9468	1
I	2655	I	São Paulo	I	8	1	4455	1
I	1741	I	Santos	I	8	1	3558	1
+-		-+-		-+-		+-		-+

## g) Qual a quantidade de clientes femininos X masculinos por cidade?

Consulta realizada:

SELECT I.cidade, c.genero, COUNT(DISTINCT c.id) AS total\_clientes FROM lojas I

LEFT JOIN transacoes t ON I.id = t.loja

LEFT JOIN clientes c ON t.cliente = c.id

WHERE c.genero IS NOT NULL AND c.genero IN ('F', 'M')

GROUP BY I.cidade, c.genero;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1 Map 2 Map 3 Reducer 4	container container	SUCCEEDED SUCCEEDED SUCCEEDED SUCCEEDED	1 1 1 1	1 1 1 1	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0
VERTICES: 04/04	[	 	>>]	100% ELAPS	ED TIME:	21.34 s		

+	+   c.genero	+	
+	+ +	++	
Agudos	I F	I 30 I	
Agudos	M	i 39 i	
Amparo	F	190	
Amparo	M	293	
Ananindeua	F	44	
Ananindeua	M	44	
Anápolis	F	46	
Anápolis	M	88	
Apiúna	F	86	
Apiúna	M	79	
Aracaju	F	1	
Araraquara	F	2	
Araraquara	M	3	
Araras	F	96	
Araras	M	86	
Araucária	F	1	
Araçatuba	F	597	
Araçatuba	M	605	
Balneário Camboriú	M	3	
Barueri	F	970	
Barueri	M	777	
Bauru	F	613	
Bauru	M	573	
Belo Horizonte	F	2679	
Belo Horizonte	M	2824	
Belém	F	424	
Belém	M	497	
Bento Gonçalves	F	257	
Bento Gonçalves	M	258	
Betim	F	276	
Betim	M	305	
Biguaçu	F	563	
Biguaçu	M	576	
Birigüi	F	82	
Birigüi	M	99	
Blumenau	F	1947	
Blumenau	M	1461	
Brasília	F	2729	
Brasília	M	2720	
Cabedelo	F	242	
Cabedelo	M	315	

#### 2. Database Netflix:

## a) Qual a média de avaliação por filme?

Consulta realizada:
SELECT f.name AS filme, AVG(a.rating) AS media\_avaliacao
FROM filme f
LEFT JOIN avaliacao a ON f.movie\_id = a.movie\_id
WHERE a.rating IS NOT NULL
GROUP BY f.name
ORDER BY media\_avaliacao DESC;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	0	 0	 0	0
Map 2	container	SUCCEEDED	5	5	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0
Reducer 4	container	SUCCEEDED	1	1	0	0	0	0
VERTICES: 04/04				100% ELAPS	ED TIME:			

```
filme
                                                        media avaliacao
| Lost: Season 1
                                                    | 4.665432098765432
                                                    | 4.589824034920202
 The Simpsons: Season 6
| Family Guy: Freakin' Sweet Collection
                                                    1 4.520766378244747
 Six Feet Under: Season 4
                                                    | 4.461601211979955
 Inu-Yasha
                                                    | 4.457773512476008
 Stargate SG-1: Season 8
                                                    | 4.456026058631922
 The Best of Friends: Vol. 4
                                                    | 4.449167996352861
                                                    | 4.43625843780135
 The West Wing: Season 3
 Lord of the Rings: The Fellowship of the Ring
                                                    | 4.43148917942777
Gilmore Girls: Season 3
                                                    | 4.428942582488959
 Firefly
                                                    | 4.410771903698188
The Simpsons: Season 3
                                                    | 4.408041422306757
Finding Nemo (Widescreen)
                                                    | 4.395903857377832
Samurai Champloo
                                                    | 4.395629238884703
The Simpsons: Treehouse of Horror
                                                    | 4.395329752312742
The Godfather
                                                    | 4.392642600755082
                                                    | 4.375887066005341
 CSI: Season 1
Buffy the Vampire Slayer: Season 6
                                                    | 4.37184159378037
 Stargate SG-1: Season 7
                                                    | 4.3612002264578225
 Friends: Season 6
                                                    | 4.359458908299912
 The Simpsons: Bart Wars
                                                    | 4.338631465517241
 Alias: Season 1
                                                    | 4.319803600654664
                                                    | 4.318832283915284
 Nip/Tuck: Season 2
 Futurama: Monster Robot Maniac Fun Collection
                                                    | 4.315998237108858
 The Sixth Sense
                                                    | 4.3114694850355555
 Angel: Season 4
                                                    | 4.310205515545406
 Farscape: The Peacekeeper Wars
                                                    | 4.304888888888889
 The Silence of the Lambs
                                                    | 4.303809453673214
 The Simpsons: Season 1
                                                    | 4.297105894161602
                                                    | 4.291043897846299
 Curb Your Enthusiasm: Season 3
 Pride and Prejudice
                                                    | 4.2896135424184765
 Sex and the City: Season 4
                                                    | 4.284771634137615
 The Best of Friends: Season 1
                                                    | 4.283065908073899
 The Blues Brothers: Extended Cut
                                                    | 4.282422123331214
 Star Trek: The Next Generation: Season 6
                                                    | 4.281840380413828
 Stargate SG-1: Season 2
                                                    | 4.267909715407262
                                                    | 4.265251032764311
 Braveheart
                                                    | 4.256434882827507
 The Best of Friends: Season 2
                                                    | 4.250601765177855
 Star Trek: The Next Generation: Season 7
 Stargate SG-1: Season 3
                                                      4.247497578301582
 Batman Begins
                                                    | 4.244315458757122
```

## b) Quais os filmes com a classificação acima de 4.5?

Consulta realizada:
SELECT f.name AS filme, AVG(a.rating) AS media\_avaliacao
FROM filme f
LEFT JOIN avaliacao a ON f.movie\_id = a.movie\_id
GROUP BY f.name
HAVING AVG(a.rating) >= 4.5

ORDER BY media\_avaliacao DESC;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container container container	SUCCEEDED SUCCEEDED SUCCEEDED SUCCEEDED SUCCEEDED	1 5 10 1	1 5 10 1 1	0 0 0 0 0	0 0 0 0 0	0 0 0 0 0	0 0 0 0 0
VERTICES: 05/05	[		>>]	100% ELAPS	ED TIME:	42.13 s		

#### Resposta:

## c) Qual a melhor avaliação por filme lançado em 2005?

Consulta realizada:

SELECT f.name AS filme, f.year AS ano\_lancamento, MAX(a.rating) AS rating FROM filme f
LEFT JOIN avaliacao a ON f.movie\_id = a.movie\_id
WHERE f.year = 2005
GROUP BY f.name, f.year
ORDER BY rating DESC;

0 0 (	) 0 ) 0
0 0 (	0
	U
0 0	0
0 0	0
0 0	0
ED TIME: 37.51 s	

## Resposta:

+	+	-++
filme	ano_lancamento	rating
+		-++
7 Seconds	2005	5
Alias: Season 4	2005	5
Coach Carter	2005	5
11:14	2005	5
Batman Begins	2005	5
Hostage	2005	5
Saving Face	2005	5
Empire Falls	2005	5
King's Ransom	2005	5
Beauty Shop	2005	5
The Ballad of Jack and Rose	2005	5
The Sandlot 2	2005	5
Dead Birds	2005	5
Pooh's Heffalump Movie	2005	5
Kicking & Screaming	2005	5
Look at Me	2005	5
Nobody Knows	2005	5
The Pacifier	2005	5
The Amityville Horror	2005	5
The Hitchhiker's Guide to the Galaxy	2005	5
The L Word: Season 2	2005	5
Unleashed	2005	5

## d) Qual o número de avaliações por filme?

Consulta realizada:
SELECT f.name AS filme, COUNT(a.rating) AS num\_avaliacoes
FROM filme f
LEFT JOIN avaliacao a ON f.movie\_id = a.movie\_id
GROUP BY f.name
HAVING COUNT(a.rating) > 0
ORDER BY num\_avaliacoes DESC;

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	1	1	 0	0	0	0
Map 5	container	SUCCEEDED	5	5	0	0	0	0
Reducer 2	container	SUCCEEDED	10	10	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0
Reducer 4	container	SUCCEEDED	1	1	0	0	0	0
VERTICES: 05/05	[			100% ELAPS	ED TIME:			

## Resposta:

<u> </u>	
filme	++   num_avaliacoes   ++
Pirates of the Caribbean: The Curse of the Black Pea	arl   117075
	102721
The Sixth Sense	102376
American Beauty	101450
Bruce Almighty	98545
The Silence of the Lambs	95053
Finding Nemo (Widescreen)	94235
The Italian Job	93886
Shrek 2	92893
Braveheart	91502
Ghost	87082
What Women Want	86756
The Last Samurai	86354
50 First Dates	85605
The Bourne Supremacy	85247
A Beautiful Mind	82347
Men in Black II	81371
The Matrix: Reloaded	79504
Speed	79476
Sleepless in Seattle	78996
Something's Gotta Give	77502
Man on Fire	77447
Kill Bill: Vol. 2	77312
Lethal Weapon	76147
Road to Perdition	74652
The Wedding Planner	74461
X2: X-Men United	73684
Signs	71405
Napoleon Dynamite	71117
Being John Malkovich	70208
The Recruit	69635
Patch Adams	69461
Eternal Sunshine of the Spotless Mind	69342
Liar Liar	69105
The Mummy	68783
Sideways	68756

## e) Qual a média de avaliação por ano de lançamento?

Consulta realizada:

SELECT f.year AS ano\_lancamento, AVG(a.rating) AS media\_avaliacao FROM filme f
LEFT JOIN avaliacao a ON f.movie\_id = a.movie\_id
WHERE a.rating IS NOT NULL
GROUP BY f.year
HAVING AVG(a.rating) IS NOT NULL
ORDER BY f.year;

MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
container	SUCCEEDED	1	1	0	0	0	0
container	SUCCEEDED	5	5	0	0	0	0
container	SUCCEEDED	1	1	0	0	0	0
container	SUCCEEDED	1	1	0	0	0	0
			100% ELAPS	ED TIME:			
	container container container 	container SUCCEEDED container SUCCEEDED container SUCCEEDED container SUCCEEDED	container SUCCEEDED 1 container SUCCEEDED 5 container SUCCEEDED 1 container SUCCEEDED 1	container SUCCEEDED 1 1 container SUCCEEDED 5 5 container SUCCEEDED 1 1 container SUCCEEDED 1 1	container SUCCEEDED 1 1 0 container SUCCEEDED 5 5 0 container SUCCEEDED 1 1 0 container SUCCEEDED 1 1 0	container SUCCEEDED 1 1 0 0 container SUCCEEDED 5 5 0 0 container SUCCEEDED 1 1 0 0	container         SUCCEEDED         1         1         0         0         0           container         SUCCEEDED         5         5         0         0         0           container         SUCCEEDED         1         1         0         0         0           container         SUCCEEDED         1         1         0         0         0

+	-+-		-+
ano_lancamento	1	media_avaliacao	1
1920	 	3.3920122887864825	- <del>-</del>
1925	-1	3.8184647302904566	١
1929	-1	3.6844629559510285	- 1
1930	-1	3.8130872204207202	- 1
1931	-1	3.75057368053477	-1
1934	-1	3.917944578961528	-1
1935	-1	3.927168416958382	-1
1936	-1	3.916519347868788	- 1
1938	-1	4.065483476132191	-1
1939	-1	4.107672595089205	- 1
1940	-1	3.8154218162278344	-1
1941	-1	3.9408571751667987	- 1
1942	-1	3.774837511606314	- 1
1943	-1	3.727823691460055	-1
1944	-1	3.7837876006213813	-1
1946	-1	3.9582454192115493	-1
1947	-1	3.765226756012923	-1
1949	-1	3.956489275047516	- 1
1951	-1	3.7986543100165036	-1
1952	-1	3.972993883126617	-1
1953	-1	3.7395563954644917	-1
1954	-1	4.082166977493708	- 1
1955	-1	3.865171413951902	- 1
1956	-1	3.7712713696207505	- 1
1957	-1	3.739489926633283	- 1
1958	-1	3.627589063794532	-1
1959	- 1	4.02670545071843	- 1
1960	П	3.721403967217009	- 1
1961	П	3.9251369757796097	- 1
1962	I	3.8925190194420964	- 1
1963	I	3.8642983664632946	١
1964	П	3.8438571973065683	١
1965	I	3.892512491894572	١
1966	I	3.5613019891500906	١
1967	I	3.6814814814814	
1968	I	3.616407931492421	I
1969		3.6984473835537663	- 1