

Análise de Dados

4.^a Aula Prática Laboratorial

Mestrado Integrado em Engenharia Informática

Ano Letivo 2019/2020

Marisa Esteves

18 de Outubro de 2019



Universidade do Minho

Plano de Aula

1. Resolução da 3.^a ficha prática laboratorial pelos alunos em grupo;
2. Correção da ficha com os alunos;
3. Finalização da resolução da 2.^a ficha prática laboratorial.

Processo ETL

Definição

O processo ETL (*Extract, Transform, Load*) é um conjunto de processos que inclui a extração de dados de fontes de informação internas e externas, podendo estar em diferentes formatos, a transformação dos dados de acordo com as necessidades da organização e, finalmente, o carregamento dos mesmos numa estrutura de dados, como por exemplo um data mart ou um data warehouse.

Processo ETL

Definição

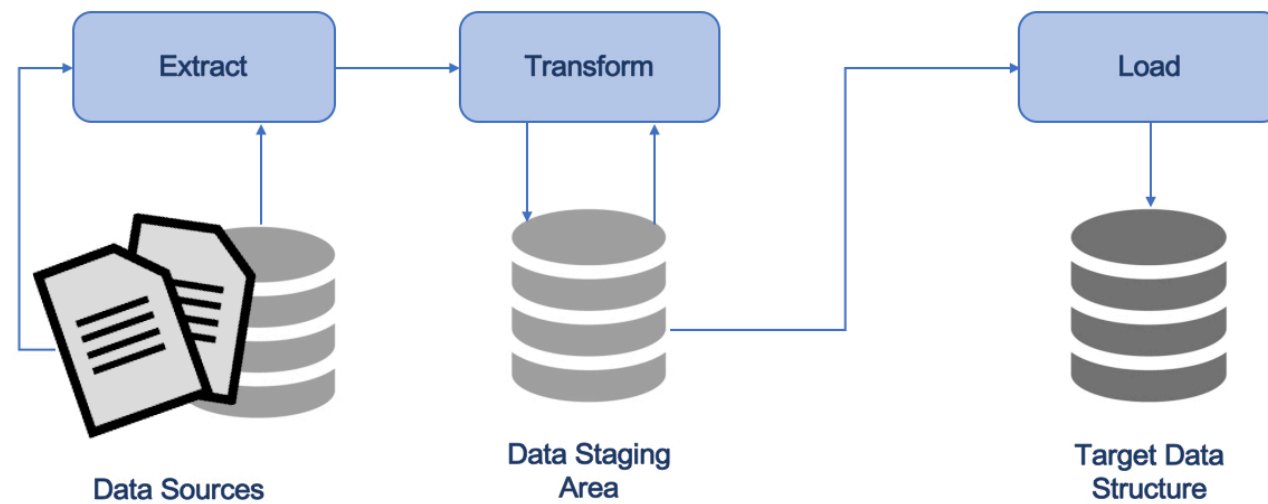


Figura 1 – Esquema do processo ETL.

Processo ETL

Porquê?

Os dados estão espalhados
por diferentes localizações

Os dados estão
armazenados em diferentes
tipos de formato

O volume de dados continua
a aumentar

Os dados podem estar
estruturados, semi-
estruturados ou não
estruturados

Data Warehousing

Definição

O processo de data warehousing enfatiza à recolha de dados de diversas fontes através do processo ETL (*Extract, Transform, Load*), correspondendo à construção de data warehouses e/ou data marts, para aceder e analisar a informação de forma útil. Os dados extraídos são processados, formatados e consolidados numa estrutura de dados única para facilitar essencialmente a análise de dados.

Data Warehousing

Definição

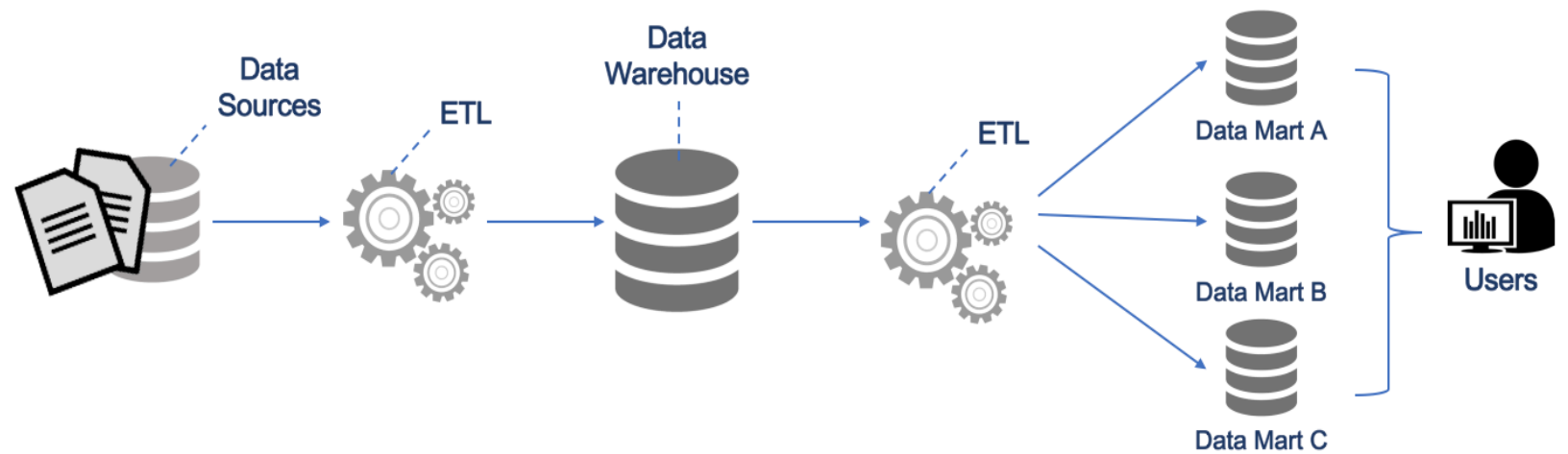


Figura 2 – Esquema do processo de data warehousing.

Data Warehousing

*Data Warehouse vs. Data
Mart*

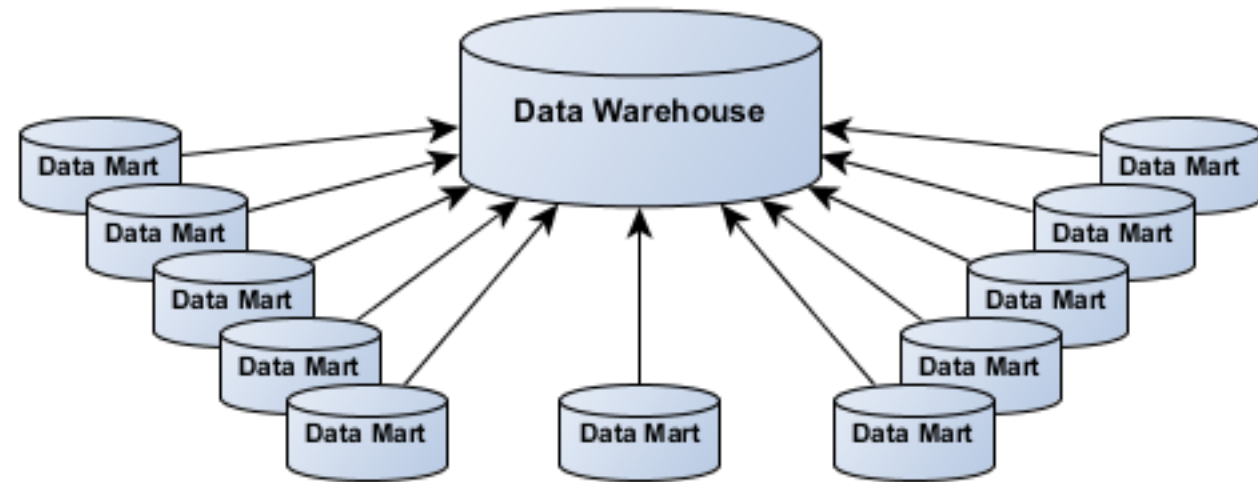


Figura 3 – Data warehouse vs. Data marts.

Data Warehousing

*Modelo Dimensional –
Esquema em Estrela vs.
Esquema em Floco de Neve*

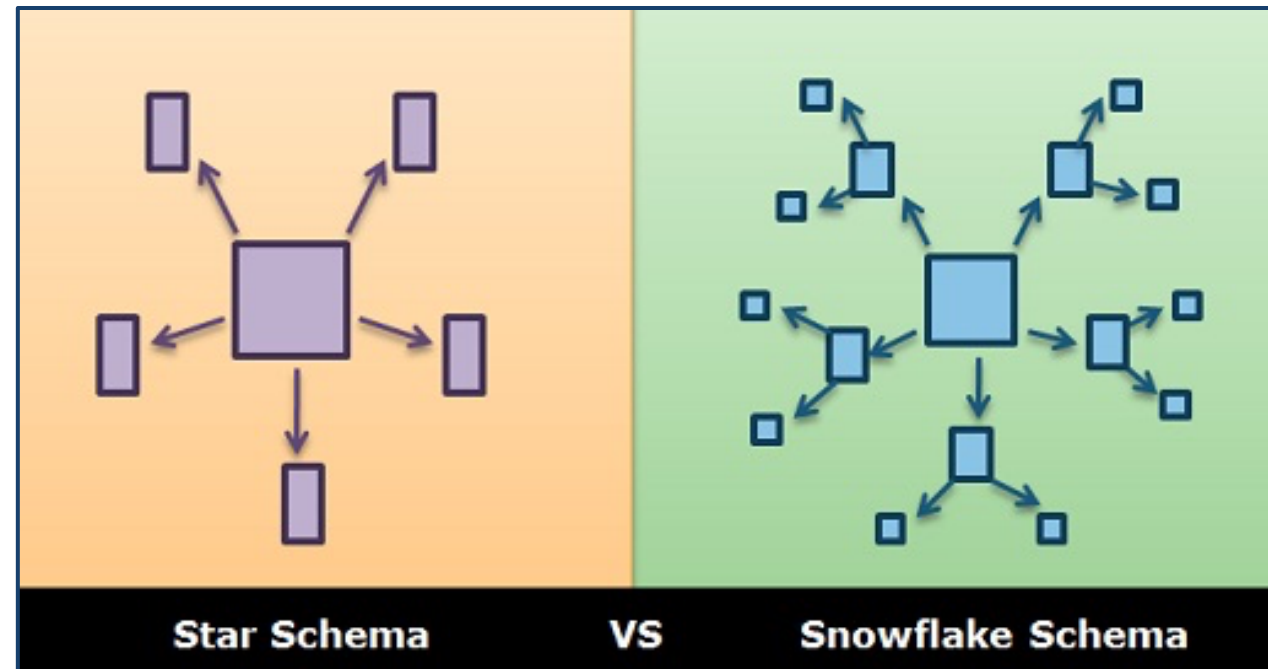


Figura 4 – Esquema em Estrela vs. Esquema em Floco de Neve.

Data Warehousing

*Modelo Dimensional –
Esquema em Constelação de Factos*

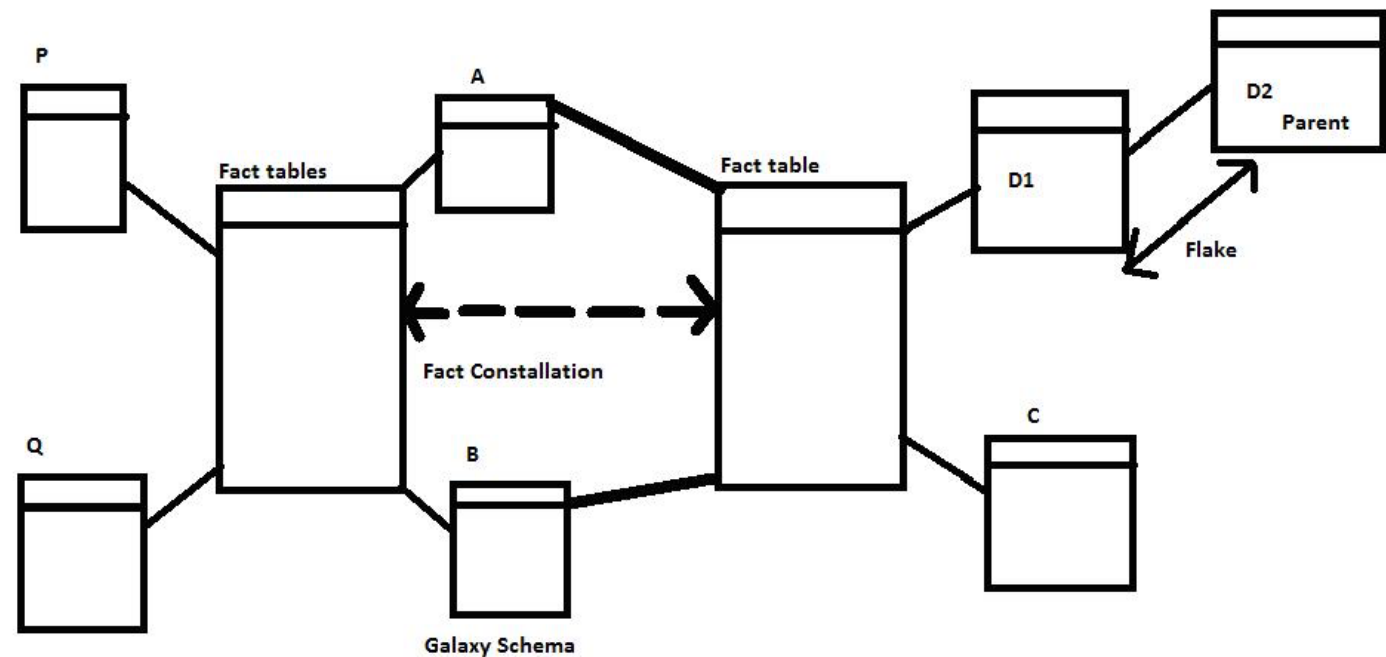


Figura 5 – Esquema em Constelação de Factos.

OLTP vs. OLAP

Definição

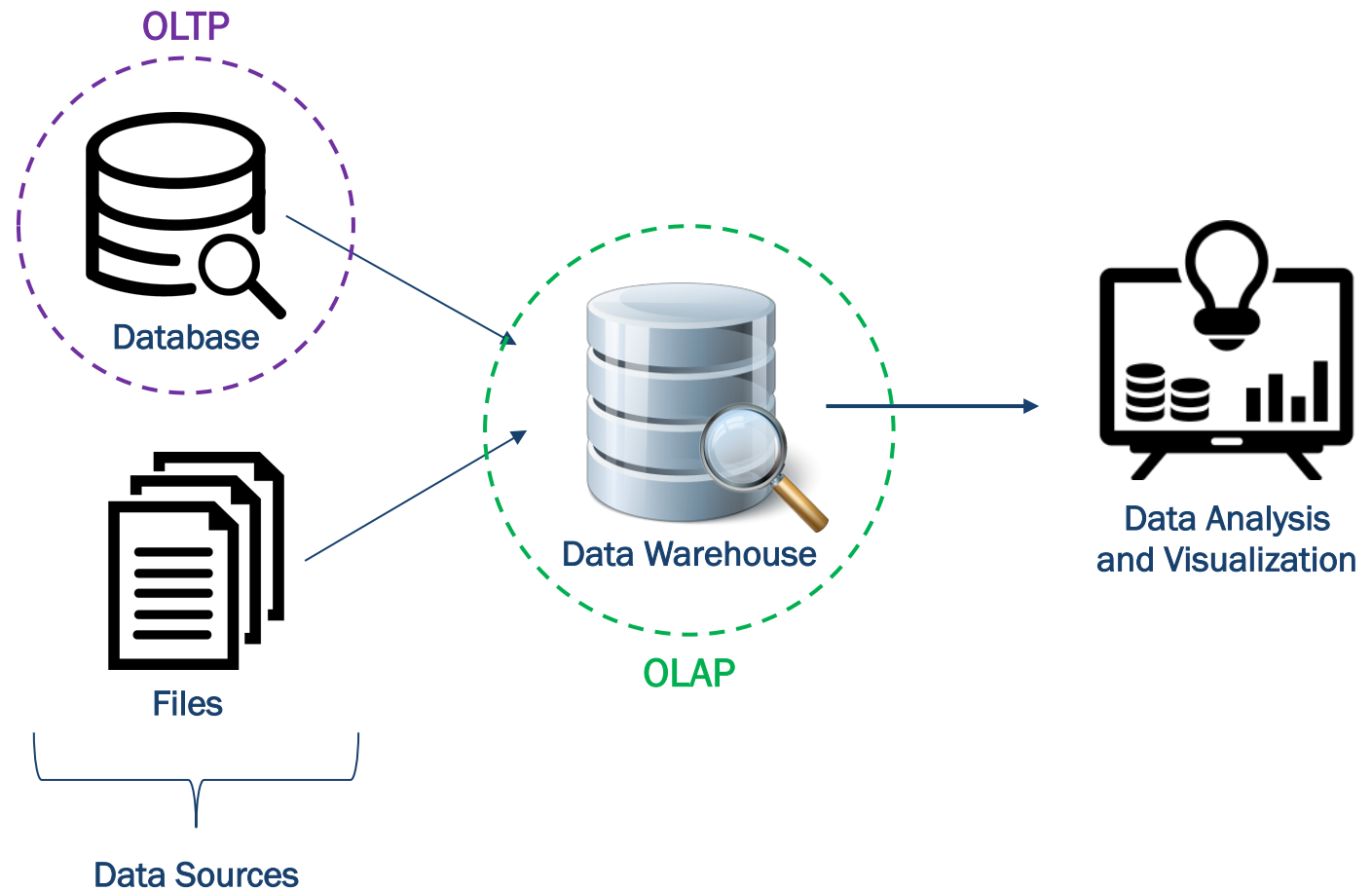


Figura 6 – OLTP (*Online Transaction Processing*) vs. OLAP (*Online Analytical Processing*).

OLTP vs. OLAP

Definição

Relational Database (OLTP)	Analytical Data Warehouse (OLAP)
Contains current data	Contains historical data
Useful in running the business	Useful in analysing the business
Based on Entity Relationship Model	Based on Star, Snowflake or Fact Constellation Schema
Provides primitive and highly detailed data	Provides summarized and consolidated data
Used for writing into the database	Used for reading data from the data warehouse
Database size ranges from 100 MB to 1 GB	Data warehouse ranges from 100 GB to 1 TB
Fast and it provides high performance	Highly flexible but it is not fast
Number of records accessed is in tens	Number of records accessed is in millions
Example: all bank transactions made by a customer	Example: bank transactions made by a customer at a particular time

Figura 7 – Diferenças entre OLTP e OLAP.

MySQL

**INSERT INTO
SELECT FROM**

Permite copiar dados de uma tabela e os inserir noutra tabela. No entanto, este comando SQL requer que os tipos de dados na tabela de origem (table1) e na tabela destino (table2) sejam iguais.

- **INSERT INTO** *table2* (*column1*, *column2*, *column3*, ...) **SELECT** *column1*, *column2*, *column3*, ... **FROM** *table1* **WHERE** *condition*

MySQL

Cursores

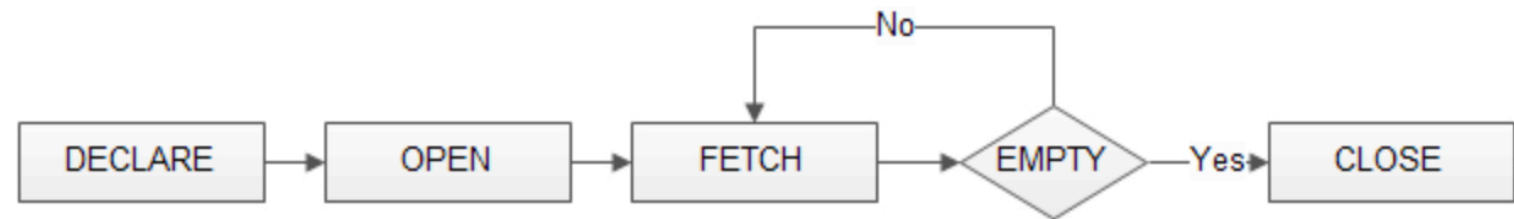


Figura 8 – Modo de funcionamento de cursores em MySQL.

MySQL

Cursores

```
1 CREATE PROCEDURE curdemo()  
2 BEGIN  
3     DECLARE done INT DEFAULT FALSE;  
4     DECLARE a CHAR(16);  
5     DECLARE b, c INT;  
6     DECLARE cur1 CURSOR FOR SELECT id,data FROM test.t1;  
7     DECLARE cur2 CURSOR FOR SELECT i FROM test.t2;  
8     DECLARE CONTINUE HANDLER FOR NOT FOUND SET done = TRUE;  
9  
10    OPEN cur1;  
11    OPEN cur2;  
12  
13    read_loop: LOOP  
14        FETCH cur1 INTO a, b;  
15        FETCH cur2 INTO c;  
16        IF done THEN  
17            LEAVE read_loop;  
18        END IF;  
19        IF b < c THEN  
20            INSERT INTO test.t3 VALUES (a,b);  
21        ELSE  
22            INSERT INTO test.t3 VALUES (a,c);  
23        END IF;  
24    END LOOP;  
25  
26    CLOSE cur1;  
27    CLOSE cur2;  
28 END;
```

Figura 9 – Exemplo de um procedimento com cursores em MySQL.

Resolução da 3.ª Ficha Prática Laboratorial

1 Esquema em Estrela vs. Esquema em Floco de Neve

O assassinato em 2014 de Michael Brown em Ferguson, Missouri, Estados Unidos da América (EUA), iniciou um movimento de protesto que culminou com o *Black Lives Matter* e um foco maior na responsabilidade dos polícias em todo o país.

Desde o dia 1 de Janeiro de 2015, *The Washington Post* tem vindo a recolher dados numa base de dados relativos a todos os disparos fatais nos EUA por um polícia durante o seu cumprimento de dever legal.

É interessante referir que é difícil encontrar dados confiáveis antes do dia 1 de Janeiro de 2015 uma vez que este tipo de acontecimentos não era documentado de forma abrangente, e estatísticas sobre a brutalidade policial estão ainda muito menos disponíveis. Como resultado, um grande número deste tipo de casos não está relatado.

The Washington Post está a recolher mais de uma dúzia de detalhes sobre cada assassinato, incluindo a raça, a idade e o género do falecido, se a pessoa estava armada e se a vítima estava num estado de crise (saúde mental).

O ficheiro disponibilizado juntamente com esta ficha prática laboratorial, nomeadamente `police_killings_us.csv`, contém os dados reais dessa recolha realizada pelo *The Washington Post*. Cada linha do ficheiro corresponde a um disparo fatal por um polícia nos EUA desde 2015. A informação representada inclui 14 colunas, nomeadamente: `id` (identificador único de cada disparo fatal), `name` (da vítima), `date` (da morte da vítima), `manner_of_death`, `armed`, `age`, `gender`, `race`, `city`, `state`, `signs_of_mental_illness`, `threat_level`, `flee` e `body_camera`.

Resolução da 3.ª Ficha Prática Laboratorial

Assim, este conjunto de dados representa uma oportunidade única para fazer questões relevantes sobre a brutalidade policial nos últimos anos nos EUA.

Com base no caso apresentado, pretende-se que:

1. Analise a estrutura da tabela `police_killings_us.csv` e, conseqüentemente, defina um modelo dimensional no formato de esquema em estrela.
2. Analise a estrutura da tabela `police_killings_us.csv` e, conseqüentemente, defina um modelo dimensional no formato de esquema em floco de neve.
3. Construa cada um dos modelos dimensionais definidos nas alíneas anteriores no MySQL Workbench (*EER diagram*).
4. Descreva as vantagens e as desvantagens entre os dois diferentes tipos de modelo dimensional definidos.
5. Defina dez questões de interesse que poderia colocar a esta base de dados. Justique a relevância de cada uma das questões definidas.