





Bloco 1

Marise Miranda



Consolidação de diferentes fontes de dados

Contexto histórico

- Produção fabril.
- Dados ad hoc.
- Dados centrados na produção.

Contexto social

- Diversidade das fontes de dados.
- Automação de serviços.
- Redes sociais.
- Conectividade e internet.

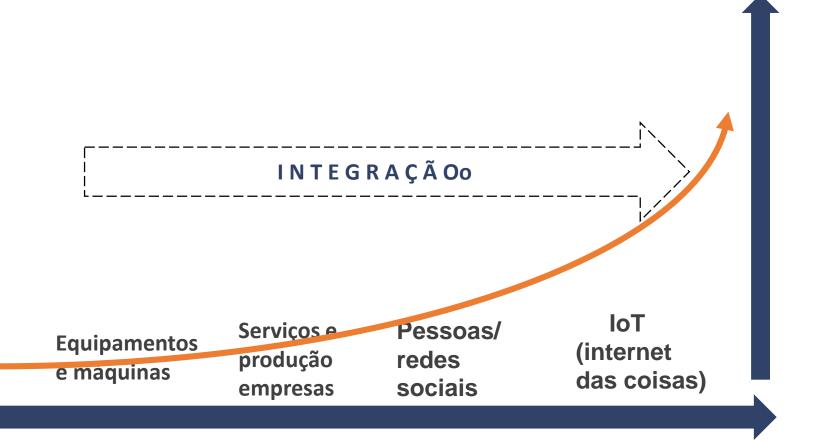
Contexto computacional

- Conectividade e internet.
- Processamento e armazenamento.



Fonte: Rawpixel/iStock.com

Dados – volume e variabilidade



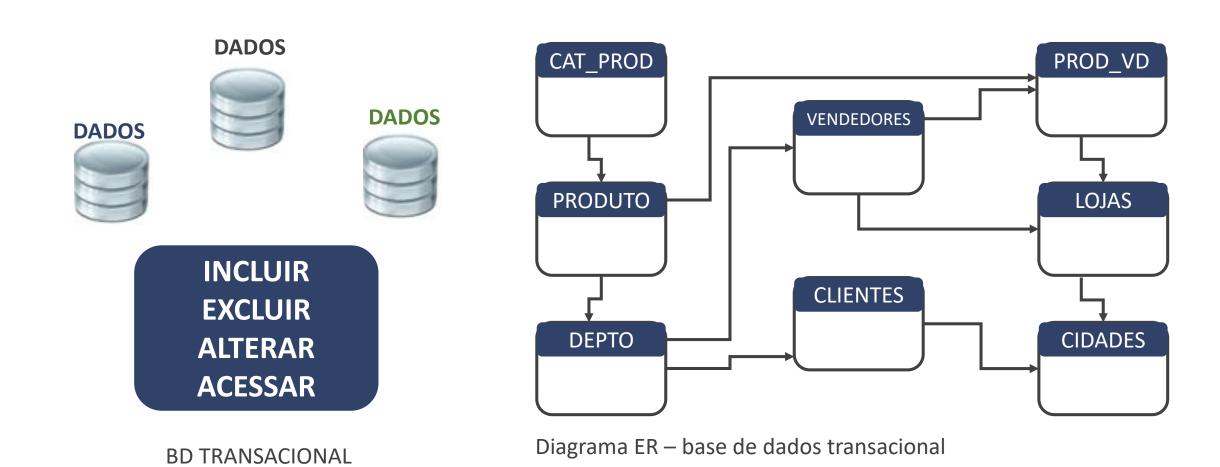
AUMENTO DO VOLUME

AUMENTO DA VARIABILIDADE

Computação, Conectividade e Internet

Fonte: elaborada pela autora

Banco de dados transacionais



▶ Banco de dados Transacionais → Modelo DW



Dimensão 1



Dimensão 2



Dimensão 3

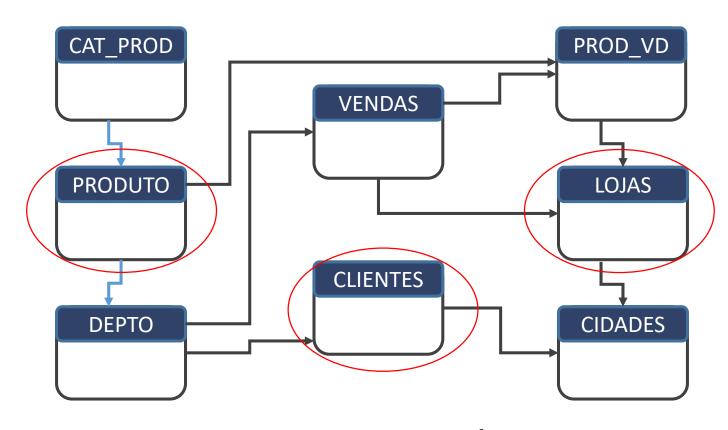


Diagrama ER − base de dados transacional → Diagrama ER do data warehouse

Modelo DW



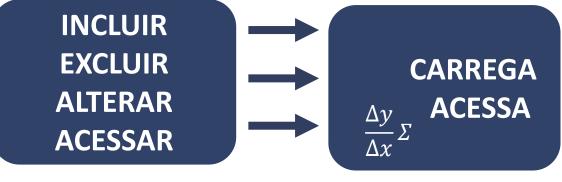
Dimensão 1



Dimensão 2



Dimensão 3



BD TRANSACIONAL

MODELO DW

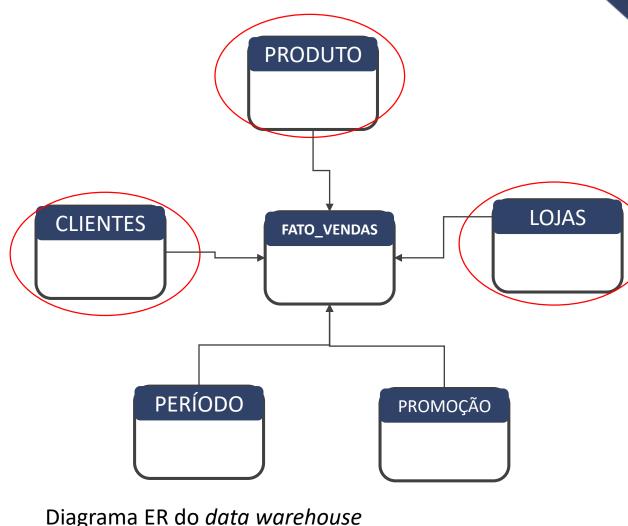
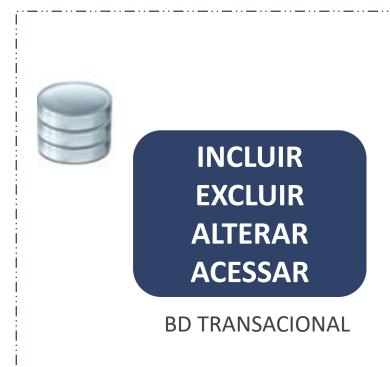
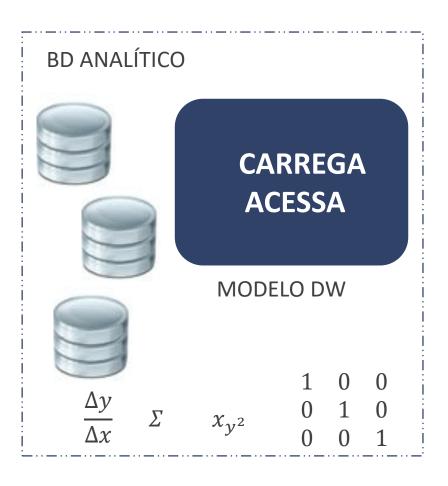


Diagrama ER do data warehouse

Modelo DW → Banco de dados analíticos





Fonte: ojogabonitoo/iStock.com



Fonte: zmicierkavabata/iStock.com

OLTP

OLAP

Conclusão

Um banco de dados transacional integra as bases de dados da organização que, em geral, advém dos processos.

Um banco de dados analítico integra um conjunto de dimensões construídas dentro da modelagem do DW. Em geral, são consumidas como apoio a tomada de decisão.

Um DW é o caminho intermediário entre os dados transacionais e os dados analíticos.



Datasets não podem ser chamados de conjunto de dados

- Datasets são a matéria prima das análises.
- Datasets são coleções de dados tabulares em formato tabela, linha e coluna, as linhas são os registros/campos e as colunas são características.
- Dataset é uma coleção de registros de dados logicamente relacionados e armazenados.
- Dataset é ad hoc, personalizado, autocontido, sem formatação.
- Datasets podem ser enriquecidos por meio de cruzamento com outros dados.
- Relatórios não são datsets.
- Conjuntos de dados tem maior abrangência.

```
46172, 0x74737769, 0x74682822, 0x72617722, 86469, 0x72656374, 0x6f72792c, 0x2073616c,
42822, 0x27257327, 0x20697320, 0x6e6f7420,
27922, 0 \times 20252064, 0 \times 69726563, 0 \times 74667279,
82073, 0 \times 65737461, 0 \times 76692070, 0 \times 66767065,
26f79, 0x616c7479, 0x0a09746f, 0x74616c20,
03d20, 0x302e300a, 0x09666f72, 0x206b2c20,
 93a0a, 0x0909746f, 0x74616c20, 0x2b3d2076,
 9202b, 0 \times 3d20765b, 0 \times 325d0a09, 0 \times 73616c65,
 0746f, 0 \times 74616c2c, 0 \times 20225673, 0 \times 65682070,
 4616c, 0x526f7961, 0x6c74795d, 0x0a092320,
06e65, 0x20646f73, 0x6567616a, 0x6f206c69,
46572, 0×6972616a, 0×20707265, 0×6b6f206b,
57620, <mark>0</mark>×0a09<mark>69</mark>66, <mark>0</mark>×206c696d, <mark>0</mark>×6974
0×696d6974,
05d20, <mark>0</mark>x3e3d206c,
                           0×65647374,
b207a, <mark>0</mark>x61207372,
```

Fonte: matejmo/iStock.com



Consolidação de diferentes fontes de dados



75,0,190,80,91,193,371,174,121,-16,13,64,-2,?,63,0,52,44,0,0,32, -0.2,0.0,6.1,-1.0,0.0,0.0,0.6,2.1,13.6,30.8,0.0,0.0,1.7,-1.0,0.6 0,5.7,-1.0,0.0,0.0,-0.1,1.2,14.1,22.5,0.0,-2.5,0.8,0.0,0.0,0.0,1 -10.0,0.0,0.0,0.6,5.9,-3.9,52.7,-0.3,0.0,15.2,-8.4,0.0,0.0,0.9,5 56,1,165,64,81,174,401,149,39,25,37,-17,31,7,53,0,48,0,0,0,24,0, ,0,0,20,0,0,0,0,0,0,0,24,52,0,0,16,0,0,0,0,0,0,0,32,52,0,0,20,0, ,0.0,7.2,0.0,0.0,0.0,0.4,1.5,17.2,26.5,0.0,0.0,5.5,0.0,0.0,0.0,0 0.9,0.0,0.0,0.4,0.7,8.3,12.3,0.2,0.0,2.2,0.0,0.0,0.0,-0.2,0.8,6. 0.9,3.8,-5.7,27.7,-0.2,0.0,9.5,-5.0,0.0,0.0,0.5,2.6,11.8,34.6,-0 54,0,172,95,138,163,386,185,102,96,34,70,66,23,75,0,40,80,0,0,24 0,0,128,0,0,0,24,0,1,0,0,0,0,0,24,36,76,0,100,0,0,0,0,0,0,0,40,2 ,0,0,0,0,0,1.0,0.0,4.5,-2.8,0.0,0.0,0.3,2.5,-2.2,19.8,0.8,-0.4,6 8.5,0.5,0.0,1.7,-2.7,0.0,0.0,-0.2,1.0,-9.4,-1.2,0.4,0.0,4.9,0.0, 0.0,5.8,-4.1,4.0,-0.5,0.4,0.3,20.4,23.3,0.7,0.0,10.0,-5.7,0.0,0. 55,0,175,94,100,202,380,179,143,28,11,-5,20,?,71,0,72,20,0,0,48, 0,60,44,0,0,32,0,0,0,0,0,56,0,0,0,0,0,0,0,0,0,0,0,0,40,44,0,0, 0,0.9,0.0,7.8,-0.7,0.0,0.0,1.1,1.9,27.3,45.1,0.1,0.0,9.1,-2.6,0. 9,0.0,3.2,-0.4,0.0,0.0,0.7,1.2,9.4,18.0,-0.1,0.0,5.1,-2.5,0.0,0. 0, -7.9, 0.0, 0.0, 0.1, 4.1, 7.6, 51.0, 0.4, 0.0, 15.0, -5.5, 0.0, 0.0, 0.1, 3.

Index of /ml/machine-learning-databases/arrhythmia

	<u>Name</u>	Last modified	<u>d</u>	<u>Size</u>	Description
	Parent Directory			_	
?	arrhythmia.data	01-Apr-1998 12	:14	393K	
?	arrhythmia.names	19-Mar-1998 16	:17	6.0K	

Apache/2,2,15 (CentOS) Server at archive.ics.uci.edu Port 443

Detalhes da Produção Intelectual Técnica de Programas de Pós-Graduação Strict...

Os dados da seção contêm informações sobre o detalhamento da produção intelectual técnica dos programas dos programas de pós-graduação stricto sensu do ano de 2017.

CSV XLSX PDF HTML

Detalhes da Produção Intelectual Bibliográfica de Programas de Pós-Graduação ...

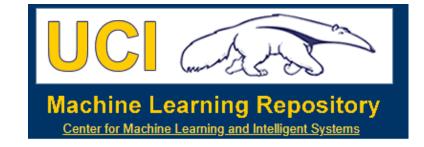
Os dados da seção contêm informações sobre o detalhamento da produção intelectual bibliográfica dos programas de pós-graduação stricto sensu no ano de 2017.

CSV XLSX PDF HTML

Discentes dos Programas de Pós-Graduação stricto sensu no Brasil 2004 a 2012

O arquivo apresenta informações sobre as atuações de pessoas declaradas como discentes pelos Programas de Pós-Graduação (PPG) stricto sensu no período de 2004 a 2012. Reúne as...

CSV XLSX PDF HTML



Conclusão

Quando se organiza um *dataset*, este deve partir do princípio do que ser quer representar. O que cada atributo significa.

Analogamente, um *dataset* deve preservar as mesmas características quando exposto à amostragem.

A abordagem por amostra dos dados não deve implicar na modelagem das análises.



Criação de um dataset - Estudo de Caso - Cafeteria

Criar dataset para análises, é uma tarefa complexa. Suponha que você tenha que criar um dataset em um Excel, relativo ao desempenho de uma cafeteria. São vendidos em média 200 cafés diariamente. A cafeteria abre de domingo a domingo, das 7 horas da manhã até às 20 horas.

Além do tradicional café, água, chá e leite também são vendidos. E ainda podem ser consumidos pão de queijo, pão com manteiga e brigadeiro.

A média de vendas de todos as bebidas, com exceção do café, é de R\$ 77,00 diariamente.

A média diária de vendas dos produtos alimentícios fica em R\$ 65,00 por dia.

Criação de um dataset – Estudo de Caso – Cafeteria

Você poderá elaborar um *dataset* simulando as vendas dos itens descritos. Leve em consideração o desempenho de um mês dessa cafeteria.

Dica: os dados podem ser preenchidos na tabela como data, hora, produto, tipo de produto, quantidade, número do pedido, ID do pedido.

Em 10 anos de vendas de café, houve uma variação de 16%, variando de R\$ 1,00 para R\$ 4,40 reais.

720.000 cafés vendidos em 10 anos, aproximadamente.

Tabule os dados e crie conceito de perspectiva analítica, siga o modelo.

Verificando as regras de negócio da cafeteria

Tabule os dados e crie conceito de perspectiva analítica. Siga o modelo.

ID pedid o	data	hor a	produt o	tipo de produto	qte	nº pedido	preço uni	preço tot
	14/01/201	07:2		café				
1	9	3	bebida	expresso	2	321	R\$4,40	R\$8,80
	14/01/201	07:2						
2	9	5	bebida	média	1	322	R\$5,40	R\$9,40
			comida	pão de queijo	1	322	R\$4,00	
3	14/01/201 9	07:2 7	comida	pão manteiga	1	323	R\$2,50	R\$2,50
4				J				
5								
6								
Etc.	Etc.	Etc.	Etc.	Etc.	Etc.	Etc.	Etc.	Etc.



Dica: 200 * 30 dias = 6000/ mês, 72.000/ano, em 10 anos = 720.000 cafés



Ferramentas SGBD de mercado e projeto Anvisa

Pesquise na internet quais são as principais ferramentas SGBDs, e compare suas vantagens e desafios. Fica como dica o Mysql e o SQL Server.





Ferramentas SGBD de mercado e projeto Anvisa

Acesse o portal da Anvisa, que apresenta o controle de medicamentos do governo brasileiro, e verifique como os dados dos produtos para saúde homologados estão disponibilizados para acesso público. Há informações sobre o projeto e quais tipos de dados estão disponibilizados, bem como o modo de como acessá-los.



Referências

INMON, W. H. Como Construir o Data Warehouse. Rio de Janeiro: Campus, 1997.

KIMBALL, R. *The Data Warehouse* Toolkit: guia completo para modelagem dimensional. Rio de Janeiro: Campus, 2002.

