

Manipulação de dados: Pandas e MySQL

Prof.: Leandro Bezerra Marinho



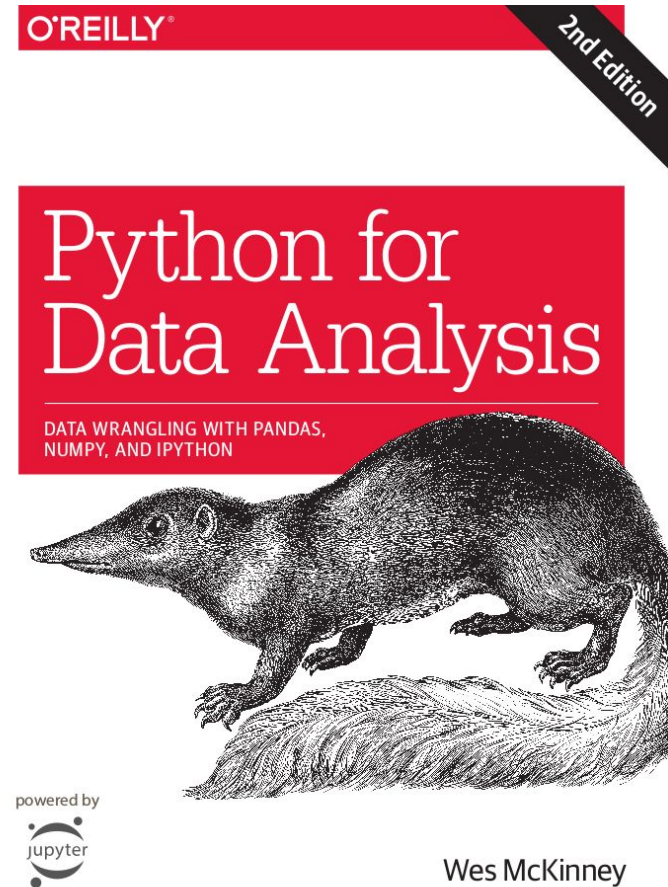
Objetivo

- Aprender os principais comandos da biblioteca Pandas.
- Entender os conceitos básicos de um banco de dados relacional.
- Compreender o SQL usando o MySQL.
- Integrar MySQL com Python.
- Conhecer as principais características do banco de dados não relacional.

O que é o Pandas?

É uma **biblioteca** open-source para **Python** que pode ser utilizada facilmente para manipulação de **dados**.

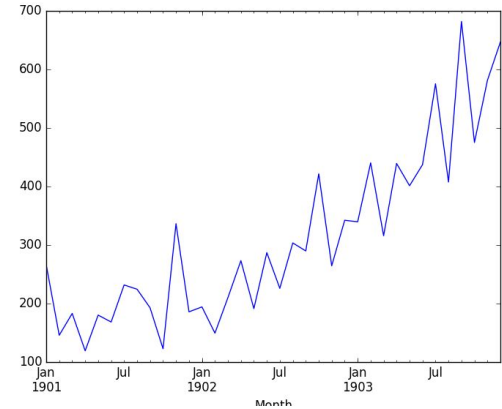
<https://pandas.pydata.org/>



Que tipo de dados?

Dados estruturados:

- Dados tabulares ou de planilhas
- Matrizes multidimensionais
- Várias tabelas de dados inter-relacionadas por colunas de chave
- Séries temporais



```
array([[[ 1., 11., 21.],  
       [nan, nan, nan],  
       [nan, nan, nan]],  
      [[ 5., 12., 22.],  
       [ 3., 13., 23.],  
       [nan, nan, nan]],  
      [[ 3., 14., 24.],  
       [ 3., 15., 25.],  
       [ 1., 16., 26.]])
```

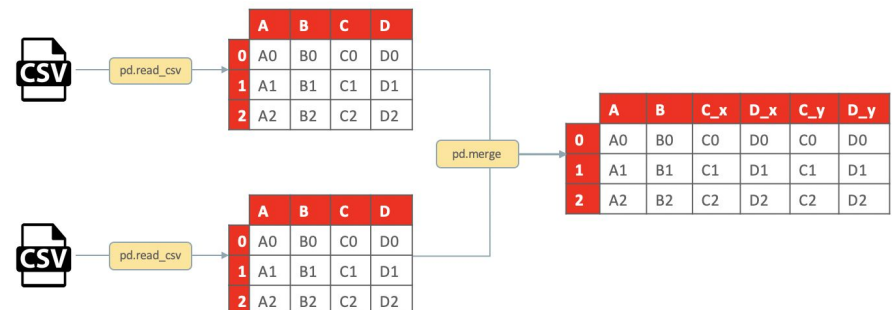
persons.csv - Notepad

File Edit Format View Help

Family Name, Given Name, VIAF ID

Ackersdijck, Willem Cornelis, 17959345	7	16	2	DM	F
Adelung, Friedrich von, 22963658	7	16	7	DM	M
Afzelius, Arvid August, 49972119	7	16	3	DM	M
Amerling, Karel, 13331054	7	16	1	DM	M
Anton, Karl Gottlob von, 183632821	7	18	3	DM	M
Arwidsson, Adolf Ivar, 8184878	7	18	7	DM	M
Asbjørnsen, Peter Christen, 116587918	7	18	4	DM	F
Attams Heinrich 37665468	7	18	4	DM	F
	7	18	7	DM	M
	7	18	7	DM	F

2013-clean 2014-clean 2014-raw



Como os dados são estruturados?

DataFrame: estrutura tabular de N-dimensões, onde cada coluna é um campo da tabela e cada linha um registro.

	<i>Name</i>	<i>Team</i>	<i>Number</i>
0	Avery Bradley	Boston Celtics	0.0
1	John Holland	Boston Celtics	30.0
2	Jonas Jerebko	Boston Celtics	8.0
3	Jordan Mickey	Boston Celtics	NaN
4	Terry Rozier	Boston Celtics	12.0
5	Jared Sullinger	Boston Celtics	7.0
6	Evan Turner	Boston Celtics	11.0

index values values values

Series: matriz unidimensional que contém uma sequência de valores acompanhado de seus respectivos índices.

	<i>Name</i>
0	Avery Bradley
1	John Holland
2	Jonas Jerebko
3	Jordan Mickey
4	Terry Rozier
5	Jared Sullinger
6	Evan Turner

index values



<https://github.com/lapisco/pandas-mysql-practices>

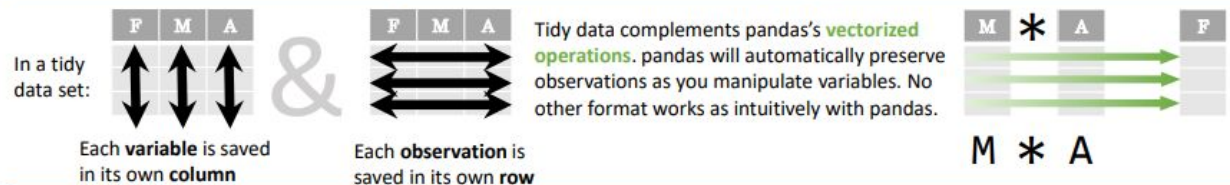
Data Wrangling

with pandas

Cheat Sheet

<http://pandas.pydata.org>

Tidy Data – A foundation for wrangling in pandas



Syntax – Creating DataFrames

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = [1, 2, 3])
```

Specify values for each column.

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
```

Specify values for each row.

n	v	a	b	c
d	1	4	7	10
e	2	5	8	11
	2	6	9	12

```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d', 1), ('d', 2), ('e', 2)],
        names=['n', 'v']))
```

Create DataFrame with a MultiIndex

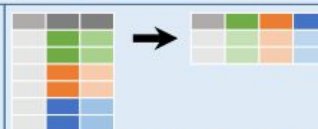
Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

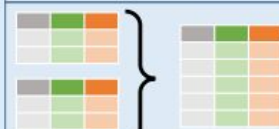
Reshaping Data – Change the layout of a data set



pd.melt(df)
Gather columns into rows.



df.pivot(columns='var', values='val')
Spread rows into columns.



pd.concat([df1, df2])
Append rows of DataFrames



pd.concat([df1, df2], axis=1)
Append columns of DataFrames

df.sort_values('mpg')

Order rows by values of a column (low to high).

df.sort_values('mpg', ascending=False)

Order rows by values of a column (high to low).

df.rename(columns = {'y': 'year'})

Rename the columns of a DataFrame

df.sort_index()

Sort the index of a DataFrame

df.reset_index()

Reset index of DataFrame to row numbers, moving index to columns.

df.drop(columns=['Length', 'Height'])

Drop columns from DataFrame

Subset Observations (Rows)



df[df.Length > 7]

Extract rows that meet logical criteria.

df.drop_duplicates()

Remove duplicate rows (only considers columns).

df.head(n)

Select first n rows.

df.tail(n)

Select last n rows.

df.sample(frac=0.5)

Randomly select fraction of rows.

df.sample(n=10)

Randomly select n rows.

df.iloc[10:20]

Select rows by position.

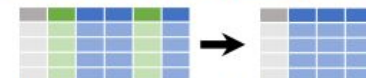
df.nlargest(n, 'value')

Select and order top n entries.

df.nsmallest(n, 'value')

Select and order bottom n entries.

Subset Variables (Columns)



df[['width', 'length', 'species']]

Select multiple columns with specific names.

df['width'] or **df.width**

Select single column with specific name.

df.filter(regex='regex')

Select columns whose name matches regular expression *regex*.

regex (Regular Expressions) Examples

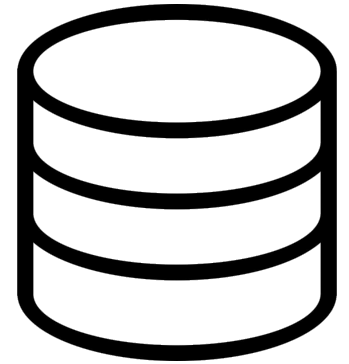
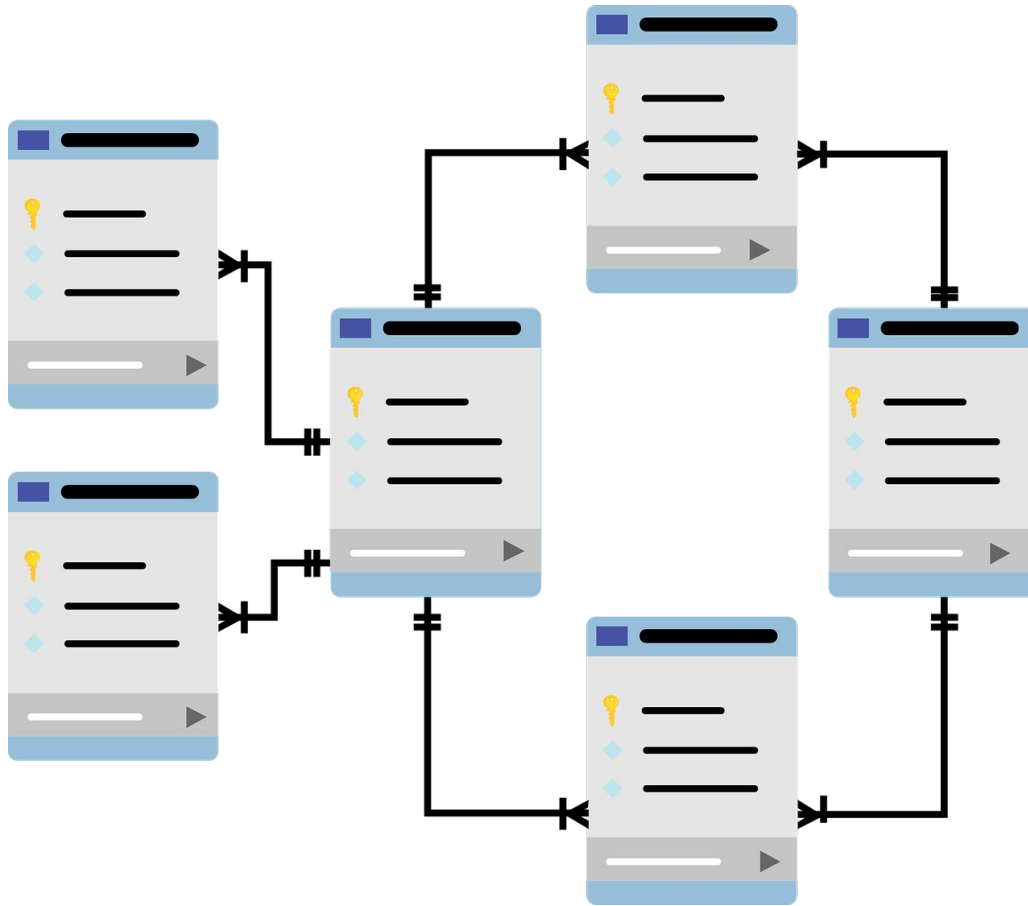
regex	Examples
'\.'	Matches strings containing a period '.'
'Length\$'	Matches strings ending with word 'Length'
'^Sepal'	Matches strings beginning with the word 'Sepal'
'^x[1-5]\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5
'^(?!Species\$).*'	Matches strings except the string 'Species'

df.loc[:, 'x2': 'x4']

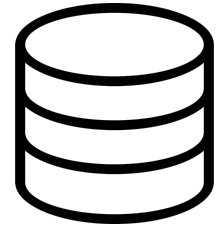
Logic in Python (and pandas)

<	Less than	!=	Not equal to
---	-----------	----	--------------

Banco de Dados Relacional

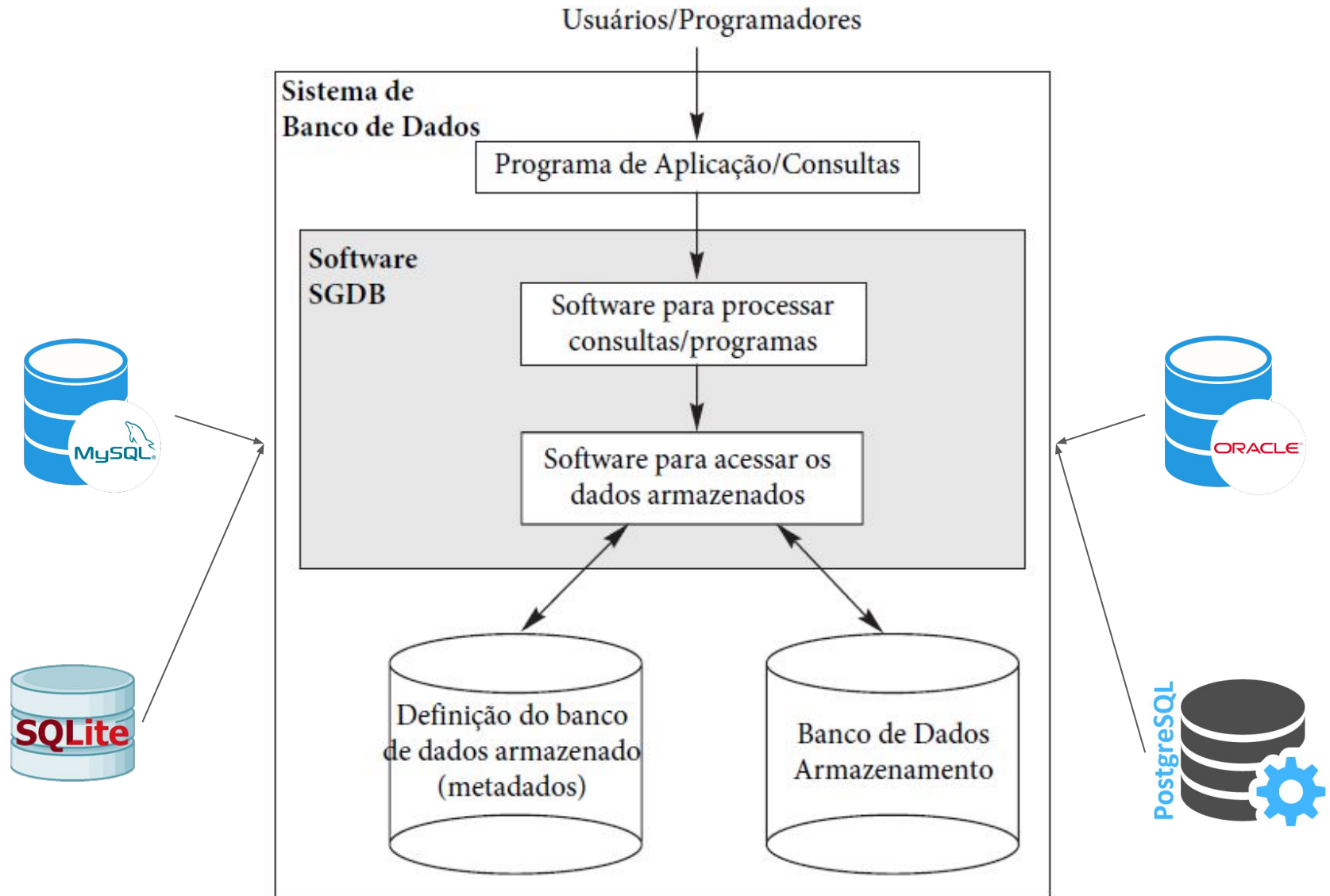


Banco de Dados Relacional



- Tipo de **banco de dados** que **armazena** e fornece acesso a dados **relacionados** entre si.
- São baseados no **modelo relacional** (tabelas).
- Cada **linha** na tabela é um **registro** com um identificador exclusivo chamada chave.
- Colunas da tabela contém **atributos** dos dados e cada registro tem um **valor** por atributo.
- Usa a Linguagem de Consulta Estruturada (SQL).

Sistema de Banco de Dados, SGBD e Banco de Dados



Vantagens do SGBD

- Estrutura de armazenamento e técnicas de pesquisa **eficientes**
- **Redundância** pode ser **reduzida**
- Dados podem ser **compartilhados**
- Suporte a **transações** pode ser fornecido
- **Backup** de recuperação
- Permite **ações** usando **regras**
- **Integridade** pode ser mantida
- **Segurança** pode ser reforçada
- **Separa** armazenamento **físico** do **conceitual**

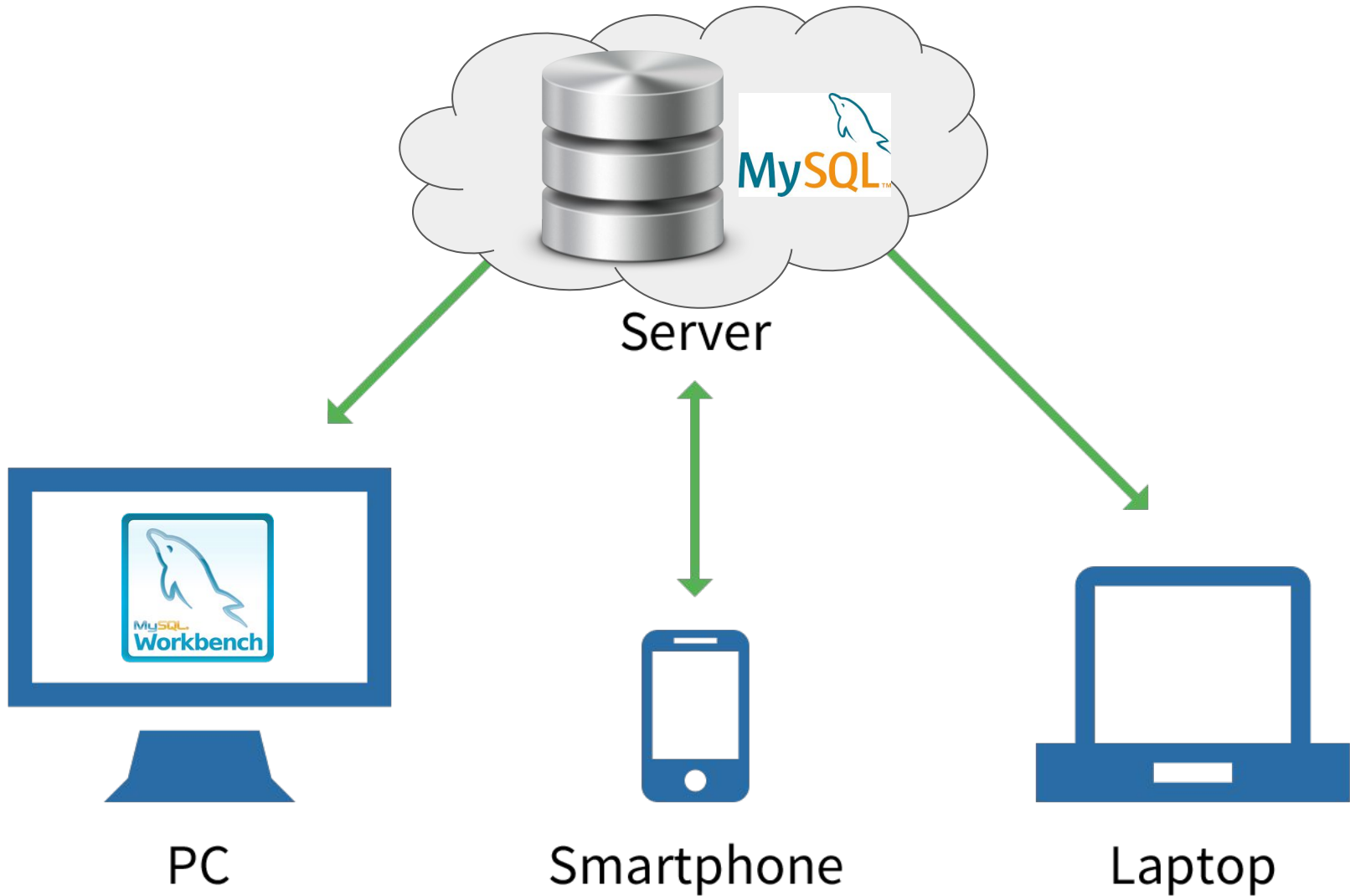


Quando não usar SGBD?

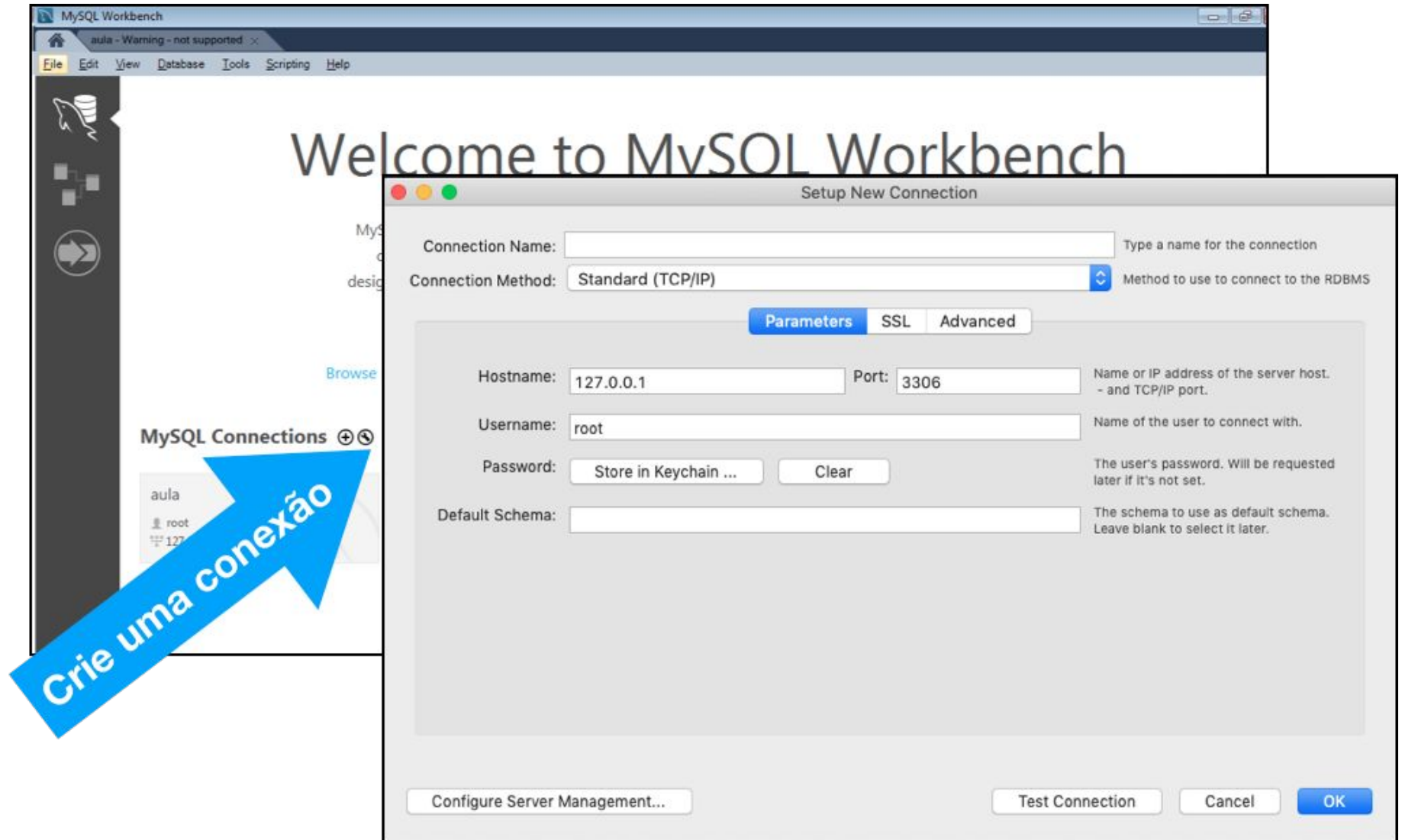
- Aplicações simples com poucas mudanças
- Sistemas embarcados com capacidade de armazenamento limitada
- Nenhum acesso de múltiplos usuários

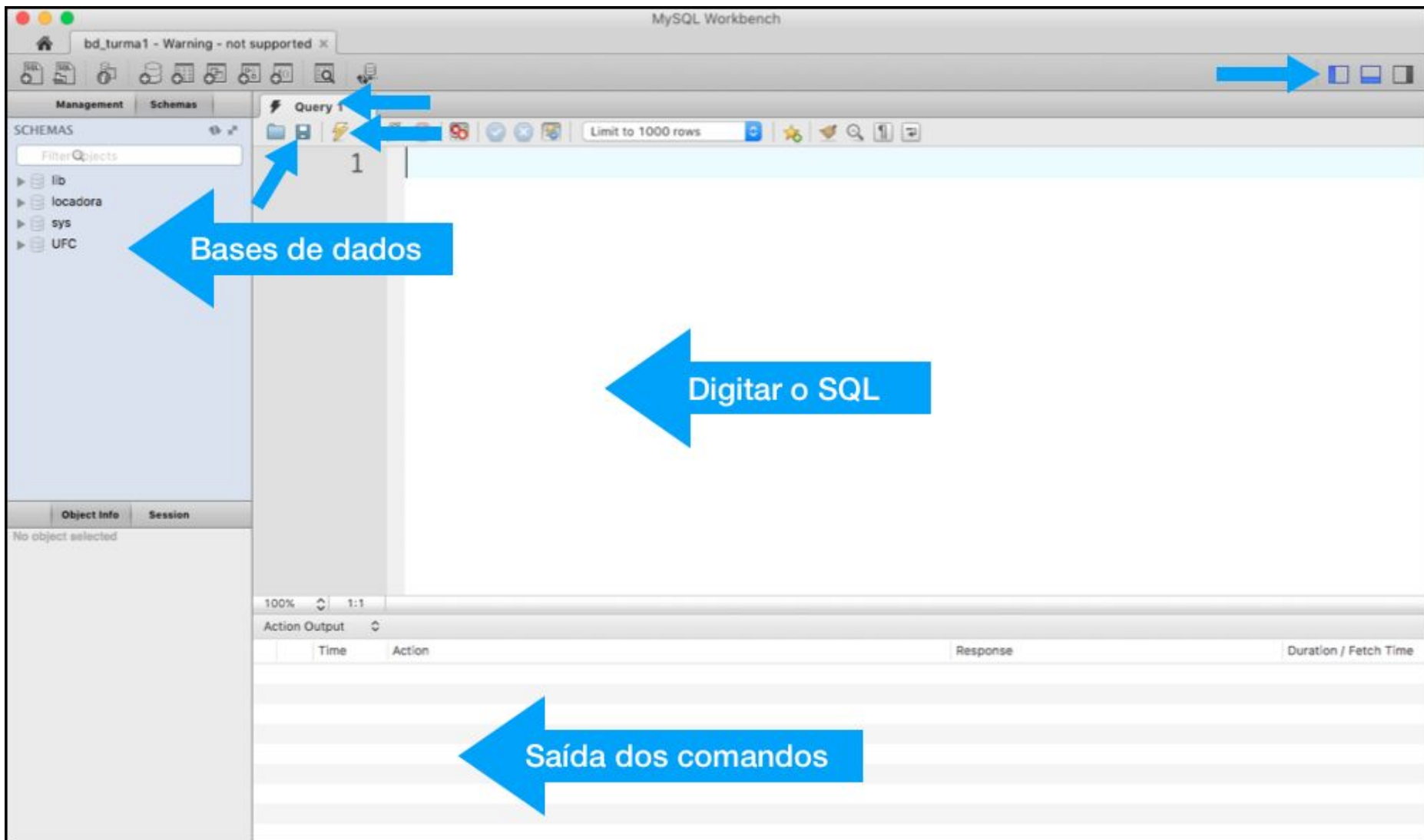


Arquitetura cliente-servidor



MySQL Workbench (cliente)







The image features a central orange rounded rectangle containing the text **.SQL** in white. A hand is shown holding the bottom of this rectangle. The background is a light gray with faint, semi-transparent SQL code snippets such as 'TABLE test (a INTEGER, b TE...', '* FROM', 'ITO to', 'ITO test (c) VALUES', and 'ITO test (b c) SEL...'. On the left and right sides, there are stylized orange database cylinders. A horizontal dotted line is positioned above the orange rectangle. At the bottom, a gray banner contains the text 'Structured Query Language'.

.SQL

Structured Query Language

Structured Query Language (Linguagem de Consulta Estruturada)

- Não é só consulta!
- Possui recursos para **definição** da estrutura de dados, **modificação** e **exclusão** de dados, etc.
- Duas principais subdivisões:

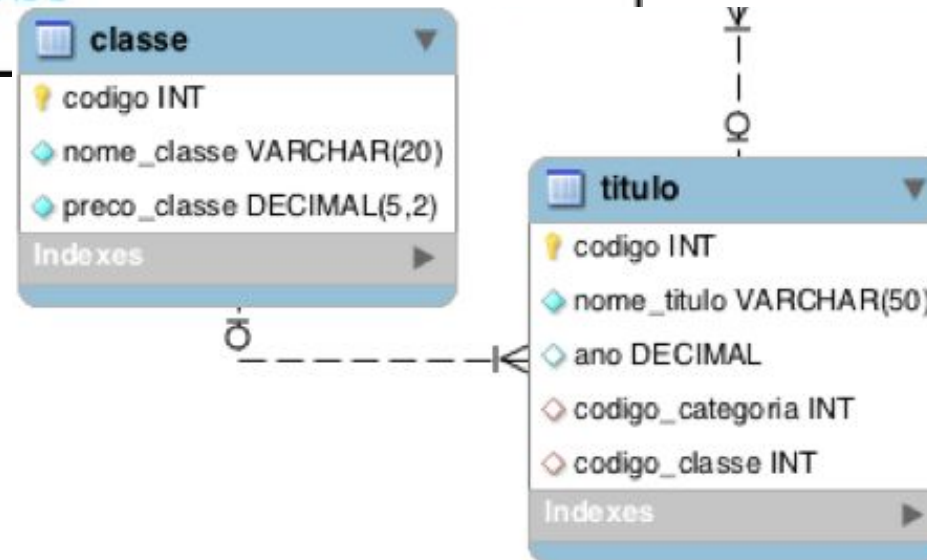
DDL
(Data Definition
Language)

DML
(Data Manipulation
Language)

Data Definition Language

```
CREATE TABLE titulo(  
  codigo          INT PRIMARY KEY,  
  nome_titulo     VARCHAR(50) NOT NULL,  
  ano             NUMERIC(4),  
  codigo_categoria INT,  
  codigo_classe   INT,  
  CONSTRAINT fk_titulo_categoria FOREIGN KEY (codigo_categoria)  
    REFERENCES categoria(codigo)  
    ON DELETE CASCADE ON UPDATE CASCADE,  
  CONSTRAINT fk_titulo_classe FOREIGN KEY (codigo)  
    REFERENCES classe(codigo)  
    ON DELETE CASCADE ON UPDATE CASCADE  
);
```

A **DDL** permite a criação (**CREATE**), remoção (**DROP**) e modificação (**ALTER**) das descrições das tabelas de uma base de dados.



Data Manipulation Language

INSERT

```
INSERT INTO projeto VALUES ('ProductA', 6, 'Umirim', 5);
```

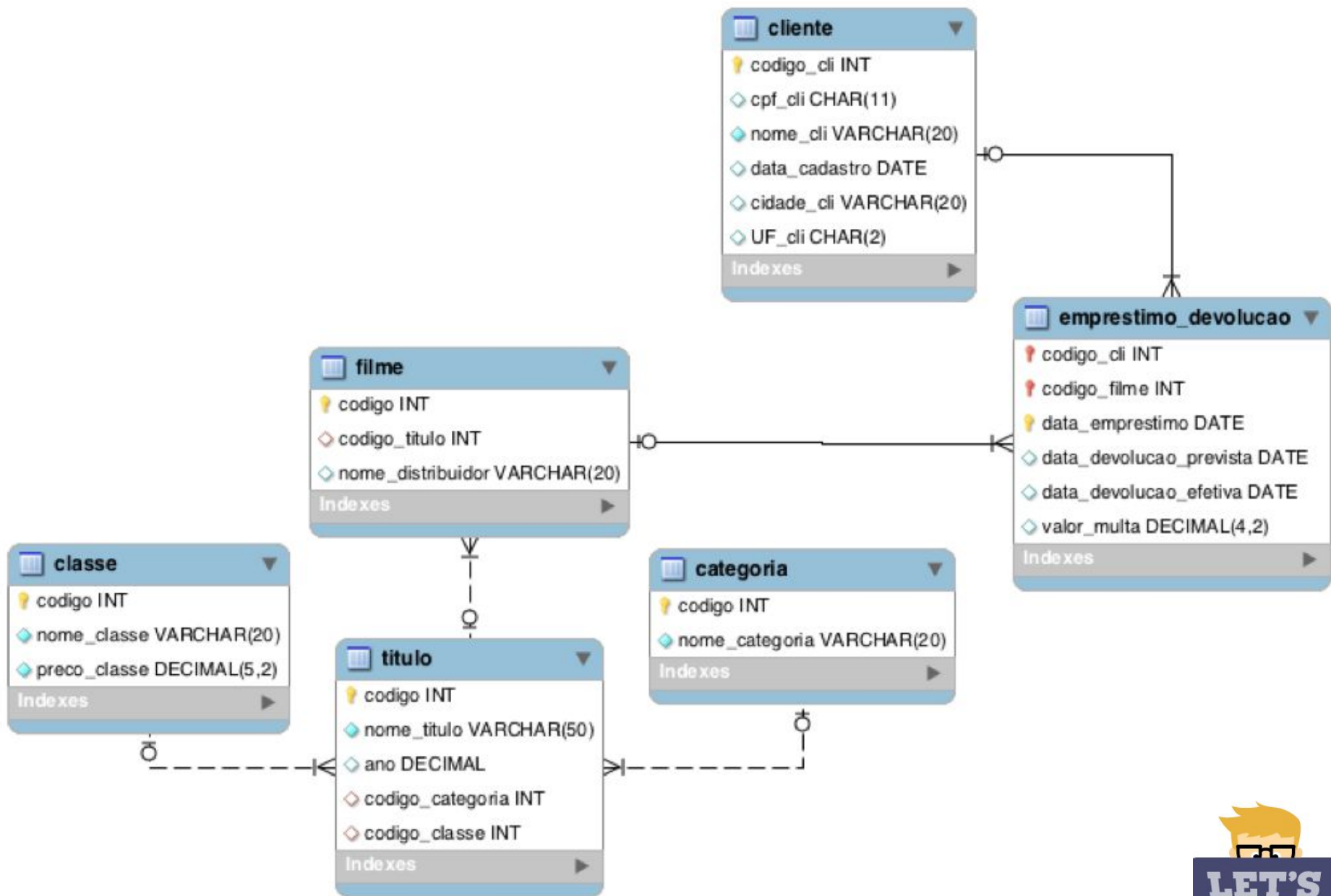
DELETE

```
DELETE FROM dependente  
WHERE nome_dependente = 'Alice';
```

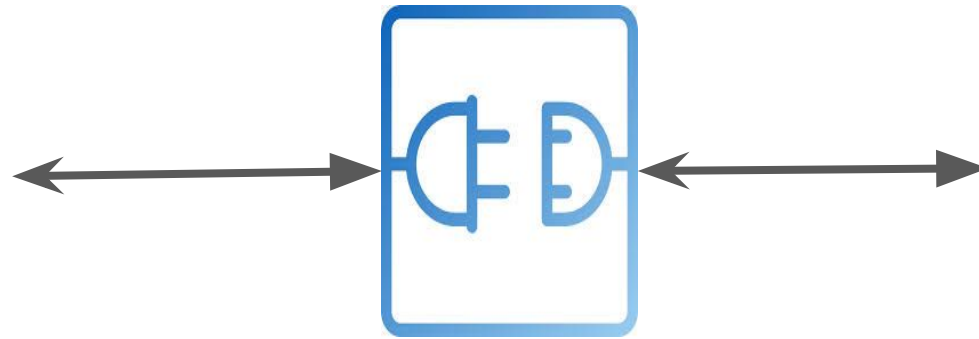
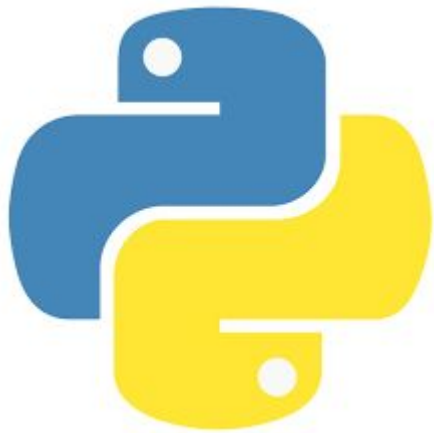
UPDATE

```
UPDATE projeto  
SET plocalizacao = 'Stafford', dnum = 5  
WHERE pnumero = 10;
```

```
31  SELECT data_nasc, endereco  
32  FROM FUNCIONARIO  
33  WHERE sexo = 'M' AND salario > 5000
```



MySQL com Python



Connector/Python





<https://github.com/lapisco/pandas-mysql-practices>

NoSQL

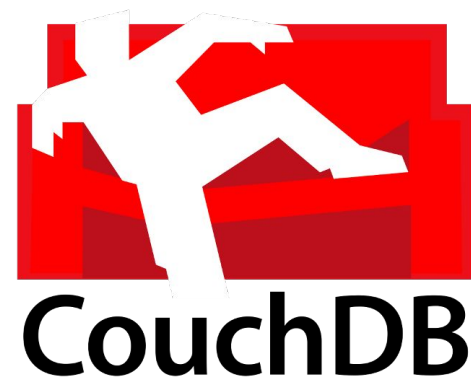


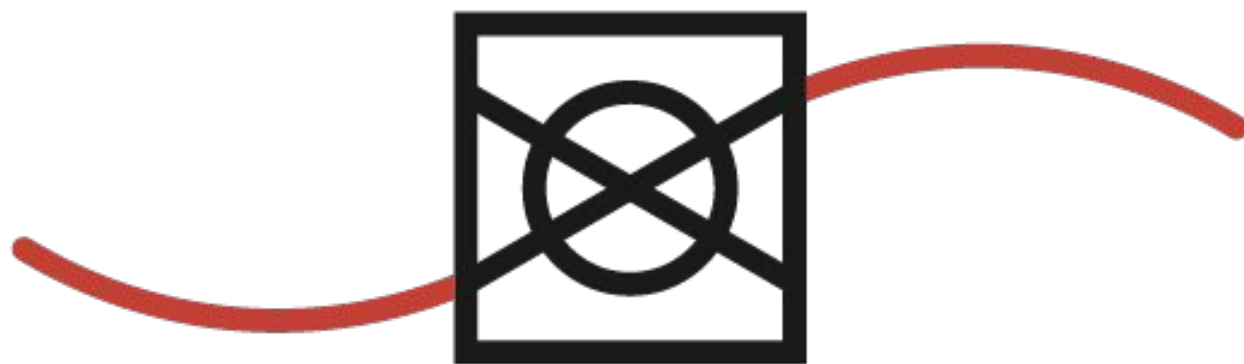
Banco de dados não relacional

- Padrão de armazenado alternativo ao modelo SQL (relacional).
- Oferecer uma maior **escalabilidade**
- Toda informação é armazenada em um registro
- Ex. do mongodb
 - `db.clientes.save({ _id: 1, fones: ["93254-8267", "93418-9592"] })`

Relacional ou não relacional, qual usar?

- NOSQL (escalabilidade)
 - Muita informação e nem todas elas são tão importantes
- Relacional (confiabilidade e a consistência)





LAPISCO

Obrigado pela atenção!