# FUNDAMENTOS E TÉCNICAS EM CIÊNCIAS DE DADOS

PROF. JOSENALDE OLIVEIRA

josenalde.oliveira@ufrn.br https://github.com/josenalde/datascience

ANÁLISE E DESENVOLVIMENTO DE SISTEMAS - UFRN



Também armazenam um KVS, mas em conjuntos permitindo o armazenamento de estruturas como um <u>JSON</u>

- 1) JavaScript Object Notation: formato compacto, padrão aberto, para troca de dados simples entre sistemas
  - 1) É completamente independente de linguagem! Douglas Crockford (1996, a partir de solução Netscape)
- 2) Cada objeto (**documento**) fica armazenado numa coleção específica, mas mesmo dentro de uma coleção não há um esquema fixo para os registros; não há **schema** (esquema)

```
"país": "Brasil",
"população": 206081432,
"PIB total em trilhões de dólares": 3.101,
"faz fronteira com": [
     "Argentina",
     "Bolívia",
     "Colômbia",
     "Guiana Francesa",
     "Guiana",
     "Paraguai",
                                 JSON
     "Peru",
     "Suriname",
     "Uruguai",
     "Venezuela"
"cidades": {
   "capital" : "Brasília",
   "mais populosa": "São Paulo"
```

Termos/conceitos do SQL	Termos/conceitos do MongoDB	
Database	Database	
Tabela	Coleção	
Linha	Documento ou documento BSON	
Coluna	Campo	
Index	Index	
Table join	Documentos aninhado ( <i>embedded</i> ) e vinculados	
Chave primária — especifica qualquer coluna única ou uma combinação de colunas como chave primária	Chave primária — No MongoDB, a chave primária é automaticamente definida como campo _id	
Agregação (group by)	Agregação de pipeline	

Relacional



NoSQL-Documento

- $\mathbf{v}$  mongo $\mathbf{DB}_{\circ}$
- \1₺ Simplicidade para transferência de informações. Tipos da dados JSON
- 2) Elemento sempre começa com chaves, e conjunto de elementos com colchetes

Tipo	Descrição	[ {	"país": "Brasil", "população": 206081432
Null	Valor vazio		
Boolean	true ou false	}, {	
Number	Número com sinal (inclui notação com E exponencial)		"país": "Argentina", "população": 41281631
String	Sequência de caracteres Unicode	}, {	
Object	Array não ordenado com itens chave-valor Chaves são strings distintas no mesmo objeto		"país": "Bolívia", "população": 10426160
Array	Lista ordenada de qualquer tipo, inteira entre colchetes e com cada elemento separado por vírgulas	}	

- 3) JSON não padroniza tipo data, nem dados binários
- 4) BSON (JSON binário variante usada no mongoDB. Ao persistir dados por mongoDB, este é o formato interno. BSON adiciona os seguintes tipos:
  - MinKey, MaxKey, Timestamp tipos utilizados internamente no MongoDB;
  - 2. BinData array de bytes para dados binários;
- 3. ObjectId identificador único de um registro do MongoDB;
- 4. Date representação de data;
- 5. Expressões regulares.

```
{
"_id" : ObjectId("57e08da696535fff4a345c67"),
"timestamp" : Timestamp(1474334118, 1),
"data" : ISODate("2016-09-20T01:16:41.720Z"
}
```

Timestamp: milissegundos desde 1.Jan.1970 (Unix Epoch) - 64 bits

```
<YYYY-mm-ddTHH:MM:ssZ>
```

UTC-3 (Brasília)



Uma nota fiscal em XML pode ser convertida para modelo relacional, mas seria decomposto em dezenas de tabelas, com chaves primárias para relacionar todas estas tabelas – custo para recuperar, processar, armazenar

(2) Num banco de dados de documentos, um XML tem sua estrut<u>ura naturalmente armazenada</u>

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<receita nome="pão" tempo_de_preparo="5 minutos" tempo_de_cozimento="1 hora">
 <titulo>Pão simples</titulo>
  <ingredientes>
    <ingrediente quantidade="3" unidade="xícaras">Farinha</ingrediente>
   <ingrediente quantidade="7" unidade="gramas">Fermento</ingrediente>
   <ingrediente quantidade="1.5" unidade="xícaras" estado="morna">Água</ingrediente>
    <ingrediente quantidade="1" unidade="colheres de chá">Sal</ingrediente>
 </ingredientes>
 <instrucoes>
                                                                           XMI
    <passo>Misture todos os ingredientes, e dissolva bem./passo>
   <passo>Cubra com um pano e deixe por uma hora em um local morno.
    <passo>Misture novamente, coloque numa bandeja e asse num forno.
 </instrucoes>
</receita>
```

Dados ESTRUTURADOS: esquema fixo, normalmente tabular (planilhas, tabelas de BD)

Não-ESTRUTURADOS: sem estrutura definida e mesmo metadados podem não ser úteis para análise (texto geral, páginas web, e-mails, postagens redes sociais, imagens, áudio, vídeo) – mais comum e que mais cresce

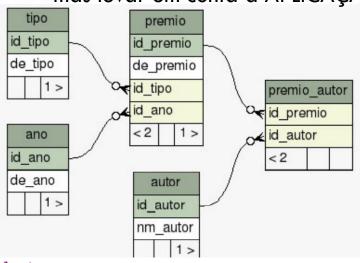
SEMI-ESTRUTURADOS: existe estrutura, mas não é fixa (XML, JSON, RDF, OWL)

```
"receita": {
 "-tempo_de_cozimento": "1 hora",
 "-tempo_de_preparo": "5 minutos",
 "-nome": "pão",
 "titulo": "Pão simples".
 "ingredientes": {
   "ingrediente": [
       "-unidade": "xicaras",
                                                  JSON
       "-quantidade": "3".
       "#text": "Farinha"
       "-unidade": "gramas",
       "-quantidade": "7",
       "#text": "Fermento"
       "-unidade": "xicaras",
       "-quantidade": "1.5",
       "-estado": "morna",
       "#text": "Água"
        "-unidade": "colheres de chá",
       "-quantidade": "1",
       "#text": "Sal"
  "instrucoes": {
   "passo": [
     "Misture todos os ingredientes, e dissolva bem.",
     "Cubra com um pano e deixe por uma hora em um local morno.",
     "Misture novamente, coloque numa bandeja e asse num forno."
```

Mais um pouco de motivação

Suponha a página da WikiPedia do <u>IgNobel</u>:, onde são identificados o <mark>ano, tipo, autor/ganhador e descrição do prêmio</mark>

No relacional, NORMALIZADO, 5 tabelas – dados organizados, mas levar em conta a APLICAÇÃO – como ocorrem as consultas



exibido d quantidad problema

#### select

and pa.id\_autor = au.id\_autor

p.de\_premio, t.de\_tipo, a.de\_ano ,au.nm\_autor
from premio p, tipo t, ano a, premio\_autor pa, autor au
where p.id\_premio = pa.id\_premio
and p.id\_tipo = t.id\_tipo
and p.id\_ano = a.id\_ano

Concordamos que um acesso ao site, se o dado que lá estivesse fosse exibido de modo relacional, pela quantidade de acessos, teria problemas de performance!

Como a página exibe tudo de uma vez, faz sentido toda a informação estar concentrada e não distribuída em tabelas! Ou seja, é preciso que os dados estejam DESNORMALIZADOS!

Aplicação obedece regras do BD, ou o BD responde às demandas da aplicação?

## mongoDB

#### Que tal esta solução?

```
"ano": 1992,
"tipo" : "Medicina",
"autores" : [
             "F. Kanda",
             "E. Yagi",
             "M. Fukuda",
             "K. Nakajima",
             "T. Ohta",
             "0. Nakata"],
"premio" :
           "Elucidação dos Componentes Químicos Responsáveis
            pelo Chulé do Pé (Elucidation of Chemical
            Compounds Responsible for Foot Malodour),
            especialmente pela conclusão de que as pessoas
            que pensam que têm chulé, têm, e as que pensam
            que não têm, não têm."
```

De 5 tabelas para 1 coleção

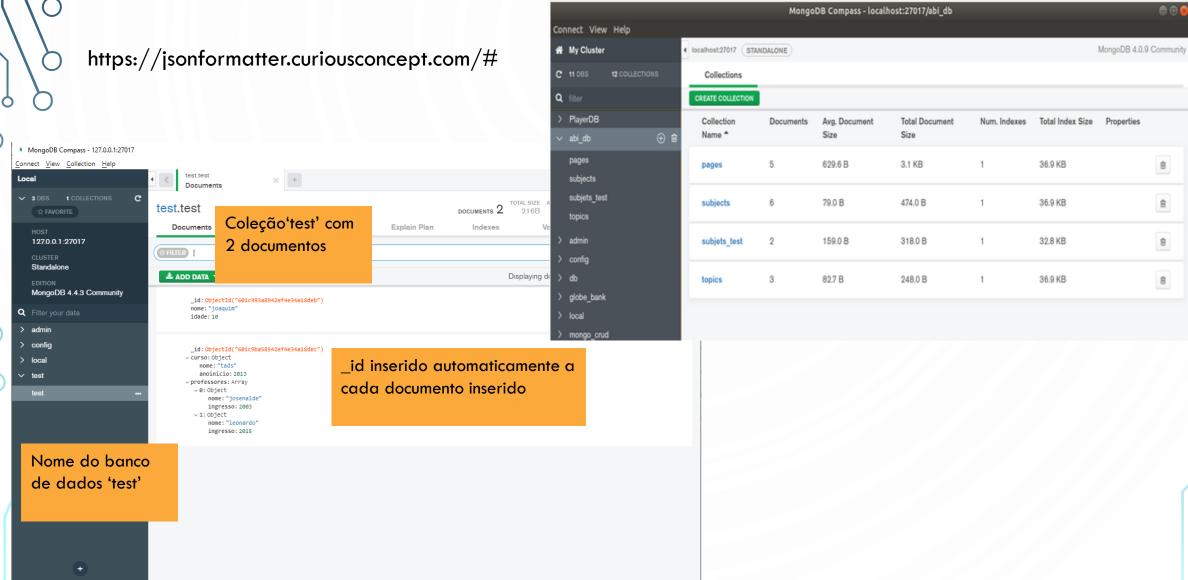
Exercício: MEGASENA

ANS LIERNI FUNDAMENTOS E TÉCNICAS EM CIÊNICIAS DE DADOS, PROF. IOSENIAI DE OLIVEIRA

MongoSH Beta



Compass



Ê

â

8

B

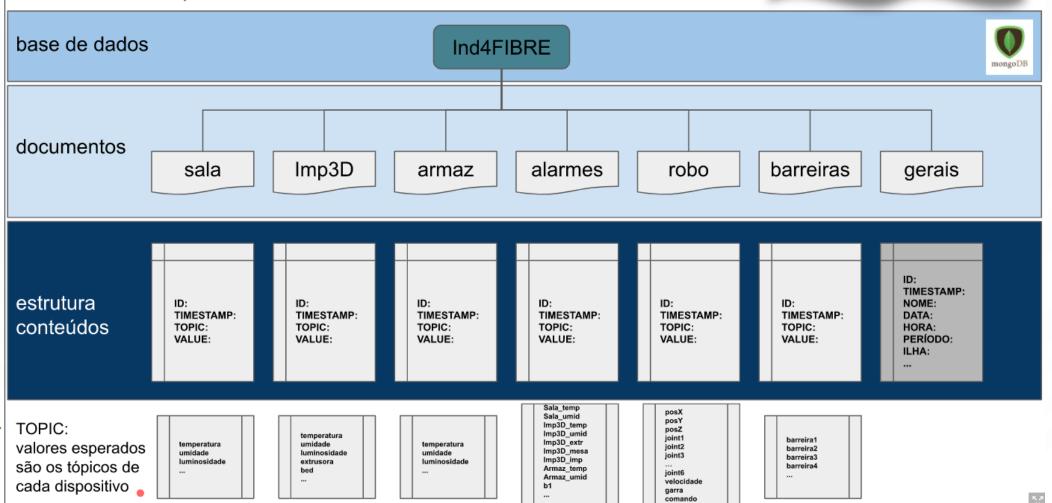
## NOSQL – DOCUMENTOS - EXEMPLO mongoDB



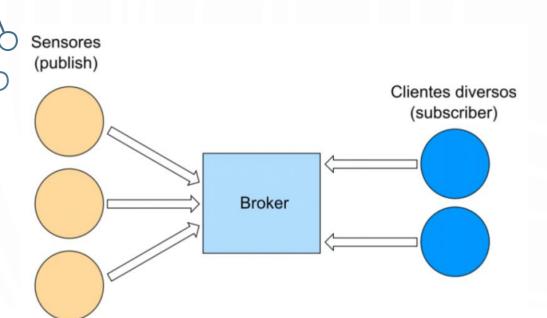
- Banco de dados local Mongodb
  - esquema

<mark>loT</mark>

Modelo orientado a documentos







Exemplo de aplicação loT: mensagens transmitidas via MQTT Mensagens podem ser encapsuladas em estrutura JSON (serialização) para transmissão e para persistência em nosql e recuperadas para exibição: https://mqtt.org/
Broker: mosquito (comum raspberry como broker)

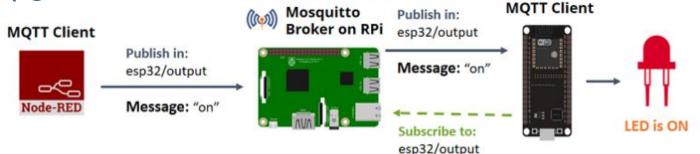
Normalmente envia msg para clientes

Pode <u>persistir a sessão do cliente</u> em memória ou em disco mosquitto.db



## NOSQL - DOCUMENTOS - EXEMPLO







- O LED está ligado ao uC em pino digital e seu estado (ON, OFF) é controlado por aplicação cliente que publica no tópico esp32/output no broker. O uC assina este tópico no broker, recebe a msg e liga o led
- 2) O uC possui sensor de temperatura e umidade conectado fisicamente e publica estas leituras nos tópicos esp32/temperatura e esp32/umidade. Estes tópicos são assinados Pela aplicação cliente, e o broker devolve ao solicitante para exibição em tela, dashboard, etc.

esp32/humidity

## >NOSQL - RESUMO



- 1) NoSQL é um termo técnico para denominar um banco de dados que não é relacional. Normalmente, ele é do tipo banco de dados de documento, orientado a objetos, chave-valor ou de grafos
- 2) De onde veio o termo mongoDB: O nome veio da palavra humongous, que significa enorme, gigantesco, para dar a ideia de grande gerenciamento de dados.
- 3) Quem usa mongoDB: Our Customers | MongoDB
- 4) Não substitui um banco relacional, pois não possui transação ou constraints de referência, que quase todo sistema possui, mas pode ser um complemento de uma base relacional, servindo como cache, por exemplo. Entretanto, se sua aplicação for desenhada adequadamente, ela pode usar inteiramente o NoSQL e não usar nenhuma base relacional.
- 5) O MongoDB cria um índice para cada collection, mas é esperado que validações nos campos aconteçam na aplicação. Existe um tipo bem simples de validação, que verifica se um campo é obrigatório ou se obedece a uma expressão regular. É possível criar índices simples e compostos (mais de um campo), inclusive para arrays. Existe também o poderoso índice de busca textual (full text search).
- 6) Opera em cluster
- O limite é de 16Mb de tamanho máximo, permitindo ter até 100 níveis de documentos aninhados. Para ter um comparativo, existem algumas versões da Bíblia em formato texto na internet, que ocupam aproximandamente 4mb. Portanto, para ultrapassar o limite atual do MongoDB, um simples registro/ documento precisa ter mais texto do que quatro bíblias completas juntas. Os nomes de campos, collections e databases podem ter até 123 bytes. Uma collection pode ter até 64 índices, cada um deles contendo entre 1 e 31 campos. O tamanho máximo do banco de dados pode variar conforme o tipo de file system e o sistema operacional. Porém, grosso modo, são 4 terabytes para Windows e 54 para Linux. Consulte os limites restantes na documentação oficial http://docs.mongodb.org/manual/reference/limits/



Construa novas aplicações com novas tecnologias





FERNANDO BOAGLIO

