



# 24º ENANCIB

Encontro Nacional de Pesquisa em Ciência da Informação  
Perspectivas Contemporâneas na Ciência da Informação

## **EXPRESSIVIDADE SEMÂNTICA DO *BIG DATA*: EXTRAINDO INFORMAÇÕES DE DADOS 24 ENANCIB, GT-2, nov. 2024, UFES, Vitória, ES**

**Grupo de Pesquisa ROCAD\* – Representação e Organização do  
Conhecimento em Ambientes Digitais, + Prof. Linair Campos,  
<http://dgp.cnpq.br/dgp/espelhogrupo/793271>**

**Carlos H. Marcondes, UFMG, UFF  
[ch\\_marcondes@id.uff.br](mailto:ch_marcondes@id.uff.br)  
<http://profmarcondes.org.br/>**

\*MARTINS, S; RAMOS JR, M; PEREIRA, D; MARCONDES, C.

# EXPRESSIVIDADE SEMÂNTICA DO *BIG DATA*: EXTRAINDO INFORMAÇÕES DE DADOS

**Questões:** o que são dados? Qual o potencial semântico dos dados? Como é gerada a semântica a partir de dados digitais?

**Hipótese:** dados se organizam em diferentes níveis, como campos, registros, tabelas, modelos conceituais, ontologias; à medida que se tornam mais complexos, em maior volume e são processados, estes agregados de dados nos diferentes níveis tornam-se mais expressivos e podem gerar semântica.

**Objetivo:** discutir teoricamente esta hipótese, propor uma conceituação da questão e mostrar indícios de sua factibilidade, propor uma aproximação CI e OC ao *Big Data* e a CD.

## Sumário

1. Introdução
2. Procedimentos Metodológicos
3. Dados e *BIG DATA* na Organização do Conhecimento e Ciência de Dados
4. Teoria dos Níveis Integrativos
5. Resultados
6. Considerações finais



# 1. Introdução, Problema

**“*While data volume proliferates, the knowledge it creates has not kept pace*”,  
Cognizant Newsletter (2011),**

**As atividades humanas são cada vez mais mediadas pelas tecnologias de informação, o chamado processo de datificação (Mejias; Couldry, 2019) da sociedade contemporânea.**

**Grande quantidade de dados digitais, chamados de Big Data, que se tornaram fundamentais nas organizações por terem alto potencial semântico.**

**Como recurso semântico, o *Big Data* só atinge seu potencial máximo quando processado por tecnologias da informação.**

# 1. Introdução, Questões de Pesquisa

- **O que são dados? Como a semântica (para humanos e máquinas) emerge de dados digitais?**

“The “strength” of the semantic, in these cases, is linked to “semantic expressivity,” associated with the tractability of the KOS for different kinds of formalism” (SOUZA et al. 2012, p. 183).

“Representational power, Semantic Expressiveness, Intelligibility” (SOUZA et al. 2012, p. 188).

“Expressivity or semantic expressivity, is a notion used to nominate how accurate a knowledge representation express a phenomena of reality”. (MARTINS et al 2024, p. 6)

- **Como a EXPRESSIVIDADE SEMÂNTICA aumenta no contexto do Big Data?**

✓ Dados são organizados em diferentes níveis, desde os dados mais simples até os agregados de dados mais complexos (Foskett, 1961; Barreto, 2008), formando conjuntos de dados ou sistemas de dados, como campos, registros, tabelas, modelos conceituais, ontologias. À medida que se tornam mais complexos e volumosos, esses agregados de dados potencialmente se tornam mais expressivos e podem gerar semântica, informações, insights para humanos e máquinas.

# 1. Introdução, definições, pressupostos

- **EXPRESSIVIDADE SEMÂNTICA** (ainda uma noção) – Precisão na representação da realidade, resultando em informação potencial, potencial informativo da SAÍDA de um processamento de conjunto de dados, imediatamente oferecida a um usuário final, permitindo-lhe ação ou tomada de decisão imediata.
- **SAÍDA** do processamento de um conjunto de dados (um agregado de dados) – “saída” de um sistema de informação imediatamente oferecida a um usuário final.

# **1. Introdução ..., Referências, Bases Conceituais**

- Web Semântica (BERNERS-LEE, HENDLER, LASSILA 2001)**
- Teoria dos Níveis Integrativos (HARTMANN 1952 ), (FEIBLEMAN 1954), (GNOLI, 2018)**
- Modos de existir SNAP e SPAN (GRENON, SMITH 2004)**
- Big Data and Knowledge Organization (IBEKWE-SANJUAN, BOWKER 2017)**
- Data, Big Data and semantics (HJØRLAND , 2018)**
- The semantic of the Semantic Web (SHETH, RAMAKRISHNAN, THOMAS 2005)**
- A taxonomy of KOS according to their semantic expressivity (SOUZA, TUDHOPE, Douglas, ALMEIDA 2012)**

# O que são dados?

Várias áreas científicas também têm feito esforços para entender, conceituar e instrumentar o Big Data, como Ciência da Computação, Ciências da Saúde e Organização do Conhecimento (Shet, 2020) (Huang et al. 2015).

**Big Data hoje X explosão de informações, o fenômeno que deu origem à CI e OC na década de 1960**

**“Big Data is the Information asset characterized by such a High Volume, Velocity and Variety to require specific Technology and Analytical Methods for its transformation into Value” Mauro, Greco e Grimaldi (2016, p. 126).**

**“Data are social artefacts” (Ibekwe-Sanjuan, Bowker 2017, p. 195)**

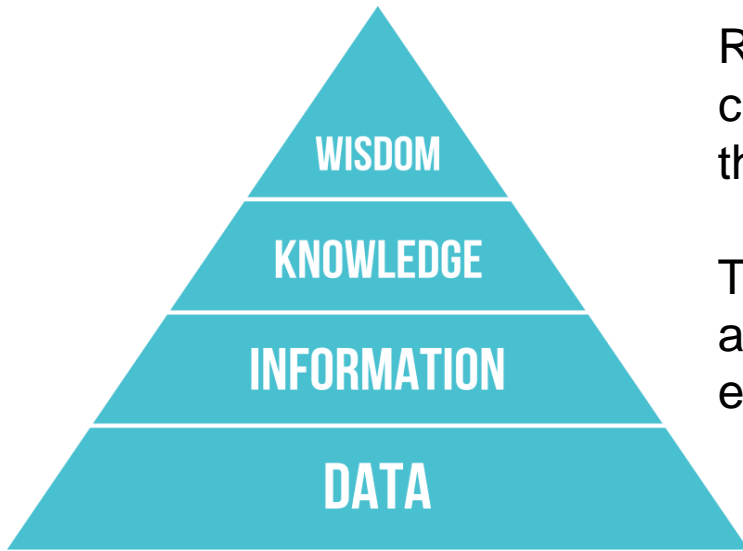
**“Dados são instanciações concretas de representações simbólicas de proposições descritivas, informadas por observação empírica, sobre as propriedades quantitativas e qualitativas de fenômenos do mundo real”; “Dados são sempre produzidos para alguns propósitos e perspectivas”. (Hjørland, 2018, s.p.).**

**“Um dado ou item de dados como um triplo  $\langle e, a, v \rangle$ , onde  $e$  é uma entidade em um modelo conceitual,  $a$  é um atributo da entidade  $e$ , e  $v$  é um valor do domínio do atributo  $a$ . Um dado afirma que a entidade  $e$  tem valor  $v$  para o atributo  $a$ ” (Hjørland, 2018, s.p.).**

**Dados são representações de entidades ou fenômenos, entidades semióticas**

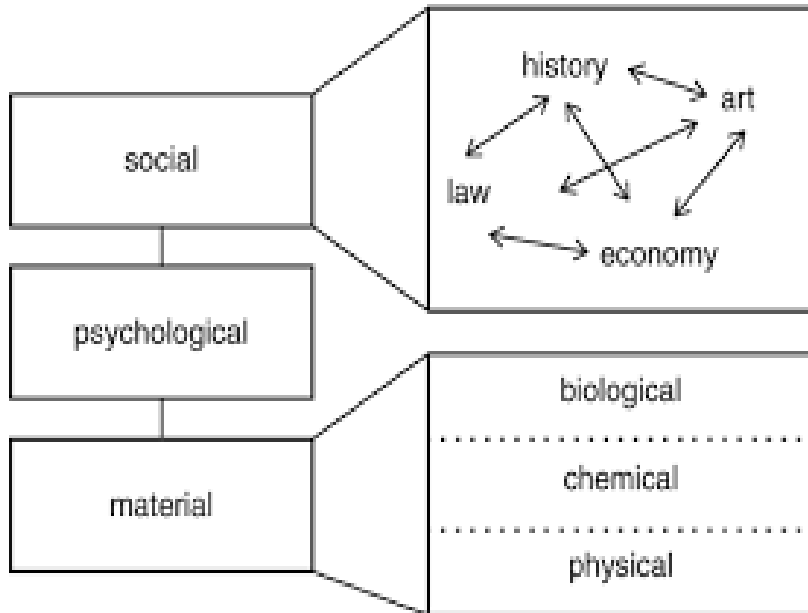


# Fundamentos ... -> Hierarquias, Teoria dos Níveis Integrativos

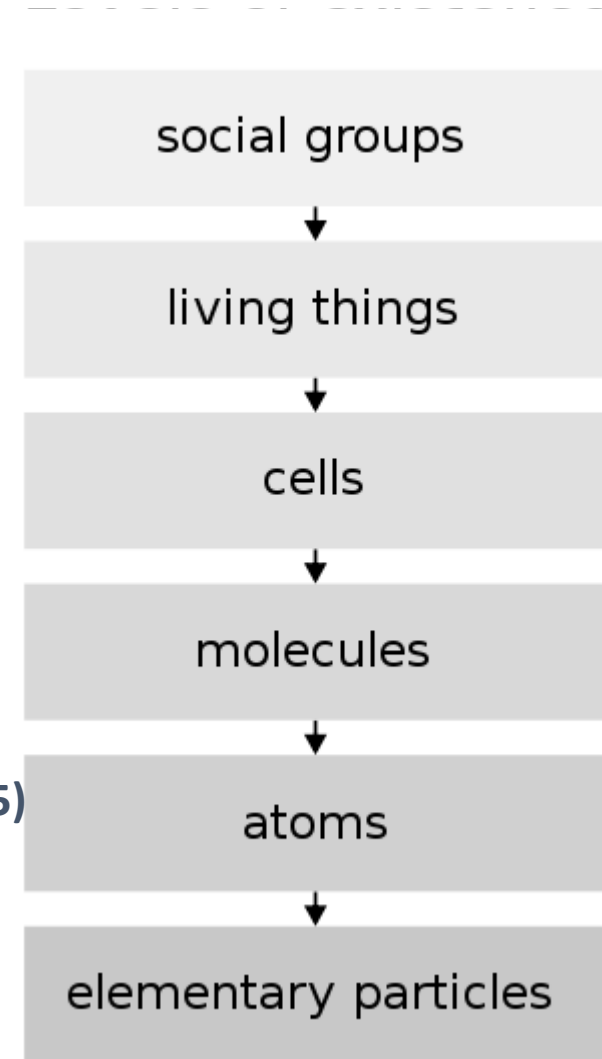


Reality is organized as crescente levels of complexity, as Phisycal level, the Biological level, the Psychologicallevel, the Cultural level

The levels overlap each other. Each higher level adds an emergent property (or quality) that did not exist at the lower level



1. Structures and Form
2. Mater and Energy
3. Cosmos and Earth
4. Life (biological systems)
5. Human beings (DAHLBERG, 1995)
6. Societies
7. Material artefacts
8. Intellectual artefacts
9. Spiritual artefacts



# Fundamentos ... -> Teoria dos Níveis Integrativos

Our list of the strata can be compared with the terminology of philosophers dealing with levels:

	Lloyd Morgan	RW Sellars	Hartmann	Poli	Popper
<b>form</b>			ideal being		
<b>matter</b>	matter	inanimate	material	material	world 1
<b>life</b>	life	animate	organic		
<b>mind</b>	mind	mind	psychic	psychological	world 2
<b>society heritage</b>		society	personal spirit objective spirit objectivated spirit	social	world 3

**DADOS**

The material and living strata can be decomposed quite easily into their layers: eg for Lloyd Morgan matter can be either physical or chemical, while mind can be conscious or reflective; for Hartmann, the spiritual stratum includes personal, objective (social), and objectivated (cultural) spirit. Modern science often acknowledges matter as including subatomic particles, atoms, molecules, celestial objects; and life as including cells, organisms, and biological populations. Layers within higher strata are less immediately identified, but levelled structures are often cited, eg families, clans, cities, nations, and the global community can be listed in the social stratum.

24 de 24- Área de transferência  
Item não coletado: exclua itens para  
aumentar o espaço disponível

# Fundamentos ...

## CONTINUANDOS E OCORENTES

GRENON, Pierre; SMITH, Barry. SNAP and SPAN: Towards dynamic spatial ontology. *Spatial cognition and computation*, v. 4, n. 1, p. 69-104, 2004

### Temporal Modes of Being

The central dichotomy among the perspectives represented in BFO concerns the modes of existence in time of the entities populating the world. BFO endorses first of all a view according to which there are entities that have continuous existence and a capacity to endure (persist self-identically) through time even while undergoing different sorts of changes. We will henceforth use the terms ‘continuant’ and ‘endurant’ interchangeably for such entities. These entities come in several kinds. Examples are: you, the planet Earth, a piece of rock; but

also: your suntan, a rabbit-hole, Leeds. All of these entities exist in full in any instant of time at which they exist at all and they preserve their identity over time through a variety of different sorts of changes. You are the same person today as you were yesterday.

In addition, however, BFO endorses a view according to which the world contains *occurrents*, more familiarly referred to as processes, events, activities, changes. **Occurrents** include: your smiling, her walking, the landing of an aircraft, the passage of a rainstorm over a forest, the rotting of fallen leaves. These entities are four-dimensional. They occur in time and they unfold themselves through a period of time.



# UM EXEMPLO...

The screenshot displays the Protege ontology editor interface. The top menu bar includes File, Edit, View, Reasoner, Tools, Refactor, Window, Ontop, and Help. The address bar shows the URL: whocovid19crfsemadatamodel (http://purl.org/vodan/whocovid19crfsemadatamodel). The main workspace is divided into several panes:

- Class hierarchy:** A tree view on the left showing the ontology structure. The class **Person** is highlighted with a red dashed box, and its subclasses are also highlighted: **WHO-COVID-19-Rapid-CRF**, **Admission form**, **Discharge/death form**, and **Follow-up**. A yellow callout box with a red border points to this hierarchy, containing the text "ESTRUTURA DO WHO-COVID-19-Rapid-CRF".
- Annotations:** A pane on the right showing the annotations for the selected class.
- Description:** A pane on the right showing the description of the selected class, including options like Equivalent To, SubClass Of, General class axioms, SubClass Of (Anonymous Ancestor), Instances, Target for Key, Disjoint With, and Disjoint Union Of.

ESTRUTURA DO  
WHO-COVID-19-  
Rapid-CRF



# NOSSO EXEMPLO:

## EXAMES

ID_PACIENTE	ATENDIMENTO	DT_COLET	DE_ORI	DE_EXAME	DE_ANALITO	DE_RESULTADO
004056afeb7b5441855846349123f686	0ae989c74676b6b5b8bf1a5f57be45f7	17/05/2020	HOSP	CULTURA AERÁ“BIA	CULTURA AERÁ“BIA	Negativa
004056afeb7b5441855846349123f686	0ae989c74676b6b5b8bf1a5f57be45f7	17/05/2020	HOSP	URÁ%IA	URÁ%IA	10
004056afeb7b5441855846349123f686	0ae989c74676b6b5b8bf1a5f57be45f7	17/05/2020	HOSP	URINA TIPO 1	pH :	7.0
004056afeb7b5441855846349123f686	0ae989c74676b6b5b8bf1a5f57be45f7	17/05/2020	HOSP	URINA TIPO 1	Urobilinogenio :	0.2

## DICIONÁRIO

	Nome	Descrição	Formato	Conteúdo, domínio e restrições
1	ID_PACIENTE	Identificação única do paciente (correlaciona com o ID_PACIENTE de todos os arquivos onde aparece (por exemplo, EXAMES e DESFECHOS)	32 caracteres alfanuméricos	string, anonimizado conforme critério da instituição
2	ID_ATENDIMENTO	Identificação única do atendimento. Correlaciona com o ID_ATENDIMENTO de todas as tabelas onde aparece (por exemplo, DESFECHO)	Alfanumérico	string, anonimizado conforme critério da instituição
3	DT_COLETA	Data em que o material foi coletado do paciente	Data (DD/MM/AAAA)	DD = Dia / MM = Mês / AAAA = Ano
4	DE_ORIGEM	Local de coleta daquele exame.	4 caracteres alfanuméricos	LAB – Exame realizado por paciente em uma unidade de atendimento laboratorial HOSP – Exame realizado por paciente dentro de uma Unidade Hospitalar
5	DE_EXAME	Descrição do exame realizado	alfanumérico	nome do exame + nome do material
6	DE_ANALITO	Descrição do analito	2 caracteres alfanuméricos	string com nome do analito
7	DE_RESULTADO	Resultado do exame, associado ao DE_ANALITO	alfanumérico	Se DE_ANALITO exige valor numérico, NNNN se inteiro ou NNNN,NNN se casas decimais Se DE_ANALITO exige qualitativo, String com domínio restrito Se DE_ANALITO por observação microscópica, String conteúdo livre
8	CD_UNIDADE	Unidade de Medida utilizada na Metodologia do laboratório específico para analisar o exame	Alfanumérico	string restrita
9	DE_VALOR_REFERENCI	Faixa de valores de referência	Alfanumérico	String - Resultado ou faixa de resultados em que é considerado normal para este analito, na população

## PACIENTES

	A	B	C	D	E	F	G
1	Column1	Column2	Column3	Column4	Column5	Column6	Column7
2	ID_PACIENTE	IC_SEXO	AA_NASCIMENTO	CD_PAIS	CD_UF	CD_MUNICIPIO	CD_CEPREDUZIDO
3	9698838a8fa8a01f15ed5c71e8e17a3	F	1962	BR	SP	SAO PAULO	CCCC
4	d9fec23b3820f93a9b1841d569db8cb5	F	1974	BR	SP	MMMM	CCCC
5	ee507ba3a9959fd31bca52852fd5715	F	1962	BR	SP	MMMM	CCCC
6	51590e8c53f4e8e332c05d7e6cee35c7	F	1960	BR	SP	SAO PAULO	CCCC
7	13699f0f7714fdaba277c5e360c6869c	M	1967	BR	SP	MMMM	CCCC
8	9bd839fe149857321685b1e1d8a55cbd	F	1948	BR	SP	SAO PAULO	CCCC
9	b723165a04c8e28fe2b4dcff8dc8dab3	F	1963	BR	SP	MMMM	CCCC

Online, real time data

Filtros

Região

UF

Município

Reg.Metropolitana/Interior

Ano

Semana Epidemiológica

Ano: 2024 X

Limpar filtros

Atualização do painel em 16/05/2024 às 15:31:58, com dados contidos nas Secretarias Estaduais de Saúde.

BRASIL

01/01/2024 a 11/05/2024

População

210.147.125

## CASOS

Casos novos notificados na semana epidemiológica

6.849

Casos Acumulados

38.802.815

Incidência covid-19 (100 mil hab)

281,68

## ÓBITOS

Óbitos novos notificados na semana epidemiológica

52

Óbitos Acumulados

712.090

Taxa mortalidade (100 mil hab)

1,64



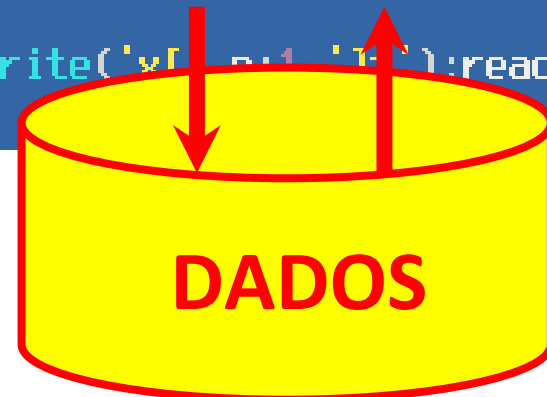
# Fundamentos ...





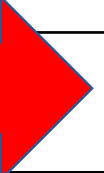



## Modelo Computacional de **Processamento de Dados**

```
Free Pascal 3.2
File Edit Search Run Compile Debug Tools
[■] SOL4.PAS
var F,L:real;
    i,j,n:integer;
    x:array[1..10] of real;
    y:array[1..10] of real;

begin
  write('n=');readln(n);
  FOR i:=1 TO n DO
  begin
    write('x[' ,i, ']=');readln(x[i]);
    write('y[' ,i, ']=');readln(y[i]);
  end;
  begin
    write('x[' ,n+1, ']=');readln(x[n+1]);
  end;
end;
```

**PROCESSAMENTO**



<b>TABELA DE NÍVEIS DE EXPRESSIVIDADE SEMÂNTICA</b>		<b>Representação</b>	<b>Processamento</b>	<b>DIMENSÃO TEMPO</b>	
<b>Concepção dos dados</b>	<b>Nível 0</b>	<b>Modelagem conceitual</b>	<b>Continuantes</b>	<b>Ocorrentes</b>	
			Processamento, consulta	Modelagem estatística	
<b>Dados digitais, arquivo digital</b>	<b>Nível 1</b>	Dados textuais	<b>PLN, REN</b> 		
Dados não estruturados	<b>Nível 2</b>	Tokens discretos, quase-signos (NÖTH, 2002), “data point” (SHAH 2020, p. 16), delimitados mas decontextualizados	<b>MIN. TEXTOS</b> 		
Dados estruturados	<b>Nível 3</b>	Dados contextualizados, um estado de coisas (JANSEN, 2008, 188), um “datum” Hjørland (2018): triplas entidade, atributo, valor, a representação de um fato ou fenômeno isolado	<b>SPARQL</b> 	<b>Ciên. Dados</b> 	
	<b>Nível 4</b>	Agregação de triplas referenciando uma única entidade ou fenômeno	<b>SQL</b> 		Dados em tempo real, um Painel
	<b>Nível 5</b>	2 ou more agregações de triplas referenciando 2 ou mais entidades ou fenômenos interrelacionados; esquema implícito	<b>SQL</b> 		
	<b>Nível 6</b>	Dados incluem o modelo conceitual/esquema; esquema explícito	<b>SPARQL</b> 		

## EXEMPLOS

**Nível 1** – “O paciente 0ae989c74676b6b5b8bf1a5f57be45f7 fez o exame URINA 1” (Um texto)

**Nível 2** – “paciente”, “0ae989c74676b6b5b8bf1a5f57be45f7”, “URINA 1”, “exame” (Tokens discretos)

**Nível 3** - <paciente\_0ae989c74676b6b5b8bf1a5f57be45f7> <exame> <URINA 1>. (Triplas RDF)

**Nível 4** – 

ID_aTENDIMENTO	DT_COLETA	DE_ORIGEM_EXAME	DE_ANALITO	DE_RESULTADO	CD_UNIDAD DE VALOR	REFERENCIA
0ae989c74676b6b5b8bf1a5f57be45f7	17/05/2020	HOSP	CULTURA AERÃ“BIA	CULTURA AERÃ“BIA	Negativa	Negativa
0ae989c74676b6b5b8bf1a5f57be45f7	17/05/2020	HOSP	URÃ‰IA	URÃ‰IA	10	mg/dL 10 a 50 mg/dL
0ae989c74676b6b5b8bf1a5f57be45f7	17/05/2020	HOSP	URINA TIPO 1	pH :	7.0	5.0 a 6.0
0ae989c74676b6b5b8bf1a5f57be45f7	17/05/2020	HOSP	URINA TIPO 1	Urobilinogenia :	0.2	mg/dl 0.2 a 1.0 mg/dl

 (Uma tabela)

**Nível 5** – 

Column1	Column2	Column3	Column4	Column5	Column6	Column7
ID_PACIENTE	IC_SEXO	AA_NASCI	CD_PAIS	CD_UF	CD_MUNICIPIO	CD_CEPREDUZIDO
9698838a8fa8a01ffd5ed5c71e8e17a3	F	1962	BR	SP	SAO PAULO	CCCC
d9fec23b3820f93a961841d569db8cb5	F	1974	BR	SP	MMMM	CCCC

 (2 ou mais tabelas interrelacionadas)

**Nível 6** –  (Uma ontologia)

**ESQUEMA**

**DADOS**

# Tipos de Saídas

- **Consulta SQL** sobre 1 ou mais tabelas interrelacionadas

**Saída** - o **subconjunto** de linhas e células de uma tabela que atendem determinada condição

- **Consulta SPARQL** em um grafo

**Saída** - o **subconjunto** das triplas que atendem determinada condição

**ESTÁTICO, Modelagem Conceitual, um campo de dados, uma célula**

**DINÂMICO, Modelagem Estatística, AO LONGO DO TEMPO, variáveis**

- **Análise descritiva** (SHAH 2020, p. 67) sobre 1 ou mais tabelas interrelacionadas

**Saída** – informações gerais sobre um conjunto de dados

- **Correlações** (SHAH 2020, p. 82) sobre 1 ou mais tabelas interrelacionadas

**Saída** – uma visão (possivelmente uma visão gráfica) sobre como as variáveis correlacionadas variam uma em relação à outra ao longo do tempo

- **Painel** (em tempo real) sobre 1 ou mais fontes de informação

**Saída** – uma visão em tempo real sobre como variáveis correlacionadas variam uma em relação a outra ao longo do tempo

# 6. Considerações finais

A semântica “emerge” dos dados em níveis crescentes de acordo com 2 eixos:

- Eixo da complexidade das representações - À medida que dados se organizam em conjuntos ou sistemas mais complexos, eles podem se tornar mais expressivos e representar coisas em um domínio com mais precisão;
- Eixo do processamento - ao processar as coisas representadas pelos dados

**MODELAGEM CONCEITUAL X MODELAGEM ESTATÍSTICA**

**Aproximar OC, RC, MC do *Big Data*, da Ciência de Dados**

## References

- LATOUR, Bruno. Redes que a razão desconhece: laboratórios, bibliotecas, coleções. In: Baratin, M; Jacob, C. O poder das bibliotecas. Rio de Janeiro : Ed. UFRJ, 2000.
- COGNIZANT. Making Sense of Big Data in the Petabyte Age. Cognizant, 20- 20 insights, jun. 2011. Disponível em: <http://www.cognizant.com/whitepapers/Making-Sense-of-Big-Data-in-thePetabyte-Age.pdf>. Acesso em: 02 abr. 2021.
- BERNERS-LEE, Tim; HENDLER, James; LASSILA, Ora. *The semantic web*. Scientific American, May, 2001.
- FEIBLEMAN, JAMES K. Theory Of Integrative Levels. The British Journal for the Philosophy of Science, v5, n.17, p. 59-66, 1954.
- GRENON, Pierre; SMITH, Barry. SNAP and SPAN: Towards dynamic spatial ontology. *Spatial cognition and computation*, v. 4, n. 1, p. 69-104, 2004.
- IBEKWE-SANJUAN, F.; BOWKER, G. C. Implications of big data for Knowledge Organization. Knowledge Organization, Baden-Baden, v. 44, n. 3, p. 187-198, 2017.
- HJØRLAND, Birger. Data with big data and database semantics. In. IEKO, ISKO Encyclopedia of Knowledge Organization. ISKO: CastanhoClaro. Disponível em: <http://www.isko.org/cyclo/data>. Acesso em: 02 dez. 2020.
- HARTMANN, NICOLAI. New ways of ontology. Henry Regnery Company, Chicago, Illinois, 1952.
- GNOLI, Claudio. Mentefacts as a missing level in theory of information science. *Journal of Documentation*, v. 74, n. 6, p. 1226-1242, 2018. Disponível em: <http://www.gnoli.eu/mentefacts.docx>. Acesso em: 07 ago. 2023.
- SHETH, Amit; RAMAKRISHNAN, Cartic; THOMAS, Christopher. Semantics for the semantic web: The implicit, the formal and the powerful. International Journal on Semantic Web and Information Systems (IJSWIS), v. 1, n. 1, p. 1-18, 2005. Disponível em: <http://www.academia.edu/download/90817267/JSWIS.pdf#page=19>. Acesso em: 12 set. 2019.
- OECD – ORGANISATION FOR ECONOMIC CO-OPERATION AND DEVELOPMENT. Handbook on Constructing Composite Indicator: Methodology and User Guide, 2008. Disponível em: <http://compositeindicators.jrc.ec.europa.eu/>. Acesso em: 22 mai. 2024.
- SOUZA, Renato Rocha; TUDHOPE, Douglas; ALMEIDA, Maurício Barcellos. Towards a taxonomy of KOo: Dimensions for classifying Knowledge Organization Systems. KO KNOWLEDGE ORGANIZATION, v. 39, n. 3, p. 179-192, 2012. Disponível em: .

**MAC Niterói**

**Museu de Arte Contemporânea, Niterói,  
Rio de Janeiro**

**Comentários são bem vindos**

**Obrigado !**

[ch\\_marcondes@id.uff.br](mailto:ch_marcondes@id.uff.br)

<http://profmarcondes.org.br>

