

# Aprendendo o ABC: Aprendizado de Máquina, Big Data e Ciência de Dados

Prof. Dr. André C. P. L. F. de Carvalho  
ICMC-USP



© André de Carvalho - ICMC/USP

1

## Tópicos

- Introdução
- Big Data
- Ciência de dados
- Aprendizado de Máquina
- Crescimento da área
- Oportunidades na área
- Ciência de dados para o bem

© André de Carvalho - ICMC/USP

2

## Presente

### Inundação de dados



<https://inside.igneous.io/storing-files-beyond-big-data>

## Causas da inundação

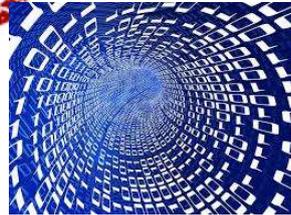
- Avanços recentes nas tecnologias para
  - **Aquisição, armazenamento, transmissão e processamento** de dados
  - Maior **quantidade**, mais **rápido**, melhor **qualidade** e menor **custo**



## Causas da inundação

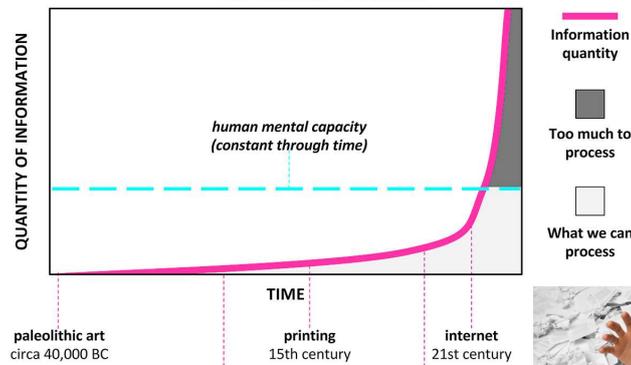
- Avanços recentes nas tecnologias para
  - **Aquisição, armazenamento, transmissão e processamento** de dados
  - Maior **quantidade**, mais **rápido**, melhor **qualidade** e menor **custo**

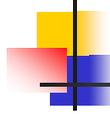
Big Data



## Sobrecarga de informação

Ryan's loose theory of Too Much Information





## Como cresce a aquisição?

Figure 1

Data is growing at a 40 percent compound annual rate, reaching nearly 45 ZB by 2020

Data in zettabytes (ZB)



Memória necessária cresce 20-40% a cada ano  
Informação duplica a cada 18-24 meses

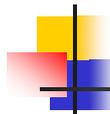
40% do crescimento se deve a atividades  
online de pessoas e empresas

**44 ZB**

Source: Oracle, 2012

André de Carvalho - ICMC/USP

7



## Como cresce a aquisição?

Figure 1

Data is growing at a 40 percent compound annual rate, reaching nearly 45 ZB by 2020

Data in zettabytes (ZB)



Em 2009, equivalente a uma pilha de  
CDs indo até a lua e voltando

Em 2020, cobrirá metade da  
distância da terra à Marte

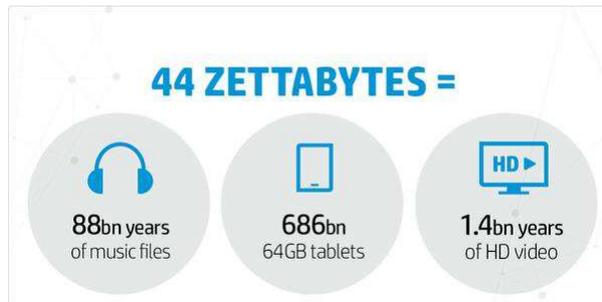
**44 ZB**

Source: Oracle, 2012

André de Carvalho - ICMC/USP

8

## O que cabe em 44 ZB?

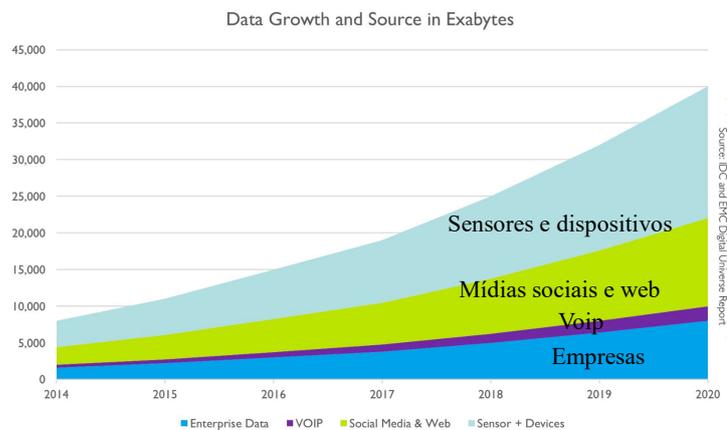


5,3 GB para cada pessoa na terra

40% desses dados precisam de proteção

<https://premioinc.com/blog/storage-overview-glimpse-digital-data-storage-technology>

## De onde vêm os dados?



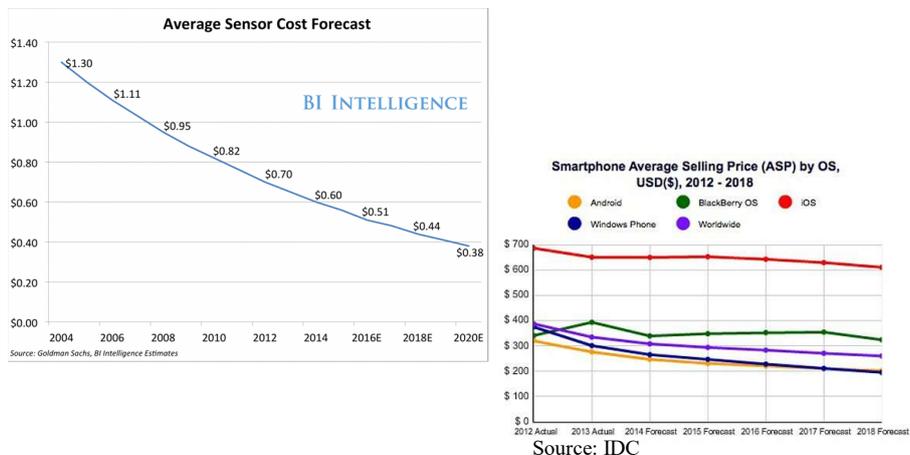
## Como segue o crescimento?



© André de Carvalho - ICMC/USP

11

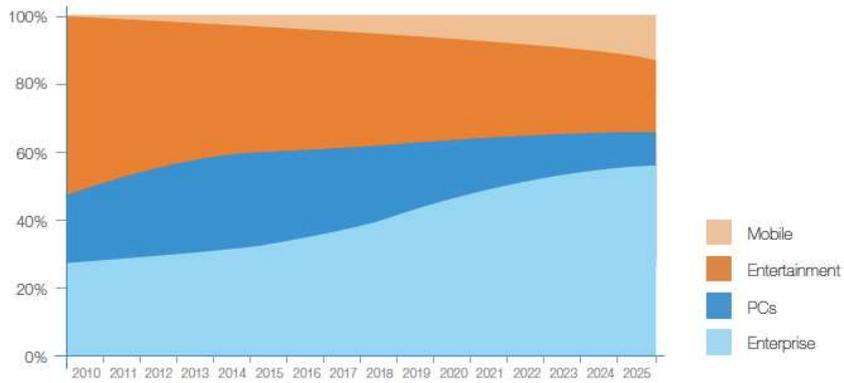
## Custo de coleta de dados



André de Carvalho - ICMC/USP

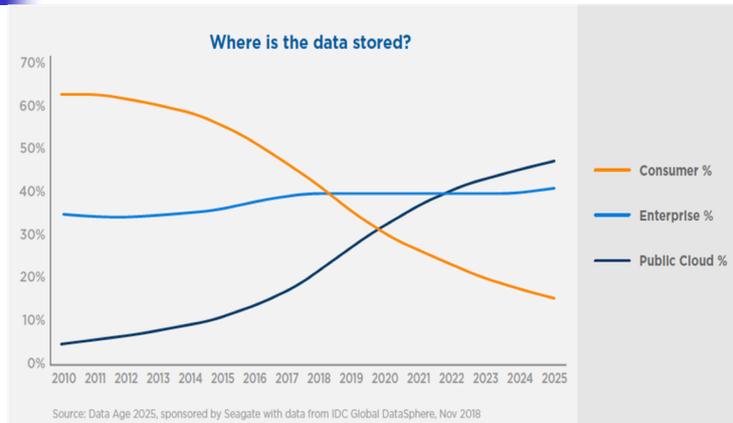
12

# Onde serão armazenados?



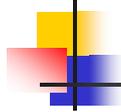
Source: IDC's Data Age 2025 study, sponsored by Seagate, April 2017

# Onde serão armazenados?

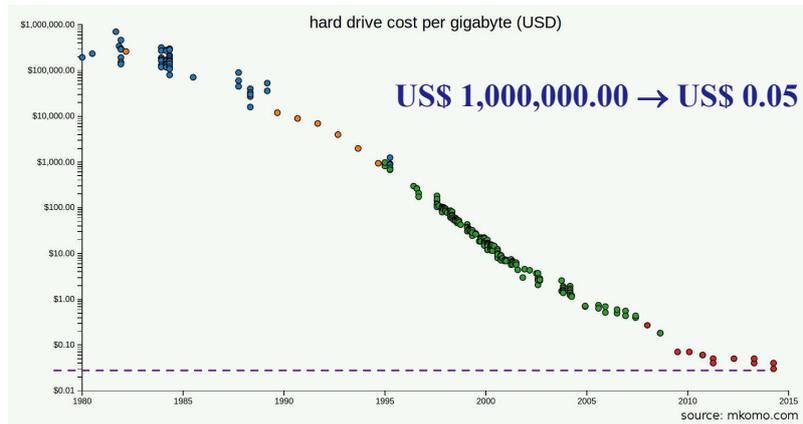


Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018

<https://www.forbes.com/sites/tomcoughlin/2018/11/27/175-zettabytes-by-2025/#736993dc5459>

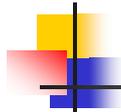


## Custo de armazenamento

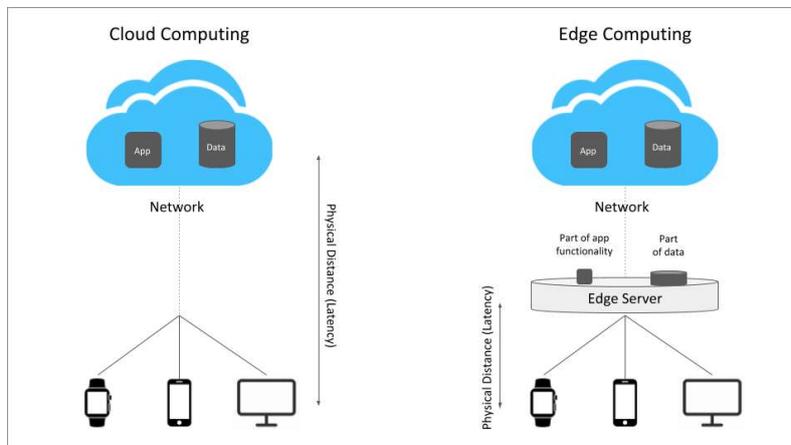


André de Carvalho - ICMC/USP

15



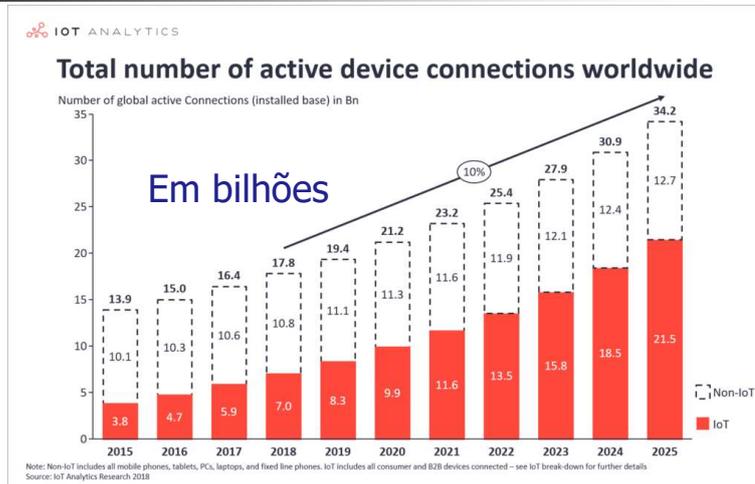
## Computação de borda



© André de Carvalho - ICMC/USP

16

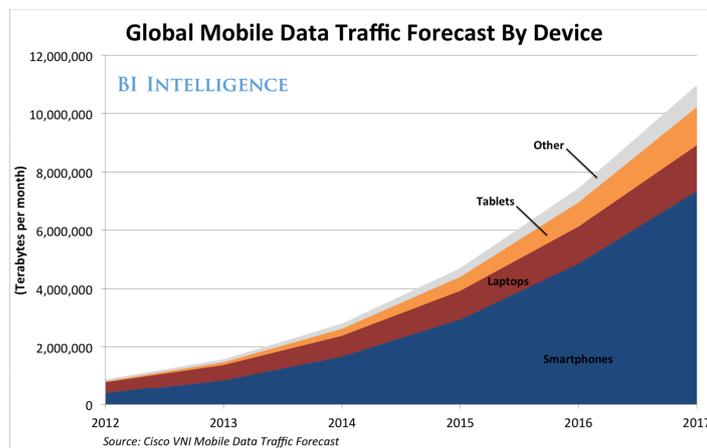
# Dispositivos conectados



© André de Carvalho - ICMC/USP

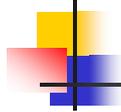
17

# Como são transmitidos?

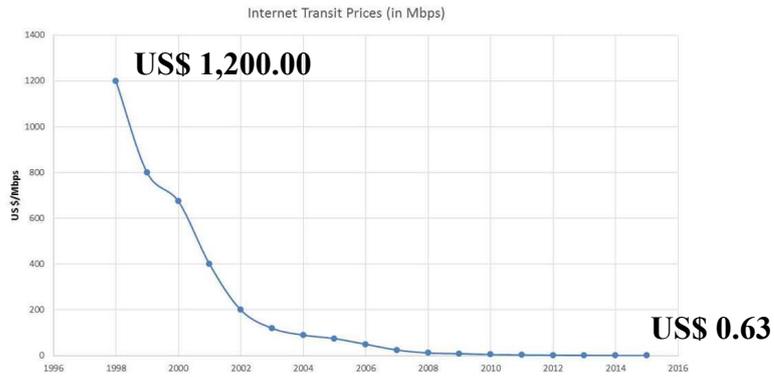


André de Carvalho - ICMC/USP

18



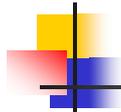
## Custo de transmissao



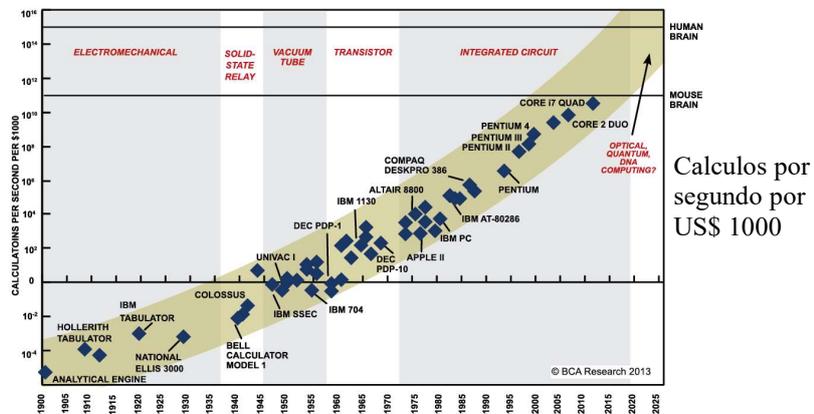
Source: DrPeering.net

André de Carvalho - ICMC/USP

19



## Capacidade de processamento



20



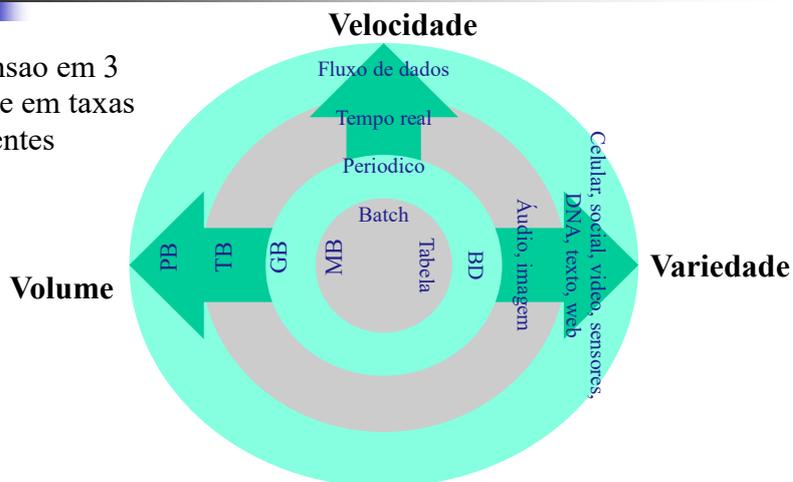
## Características de Big Data

- Grande **volume** de dados, gerados a uma grande **velocidade** e com uma grande **variedade** (3 Vs)
  - Volume: tanto de dados estruturados quanto de não estruturados
  - Variedade: vindos de fontes diferentes e com diferentes formatos
  - Velocidade: gerados em fluxos cada vez mais rápidos

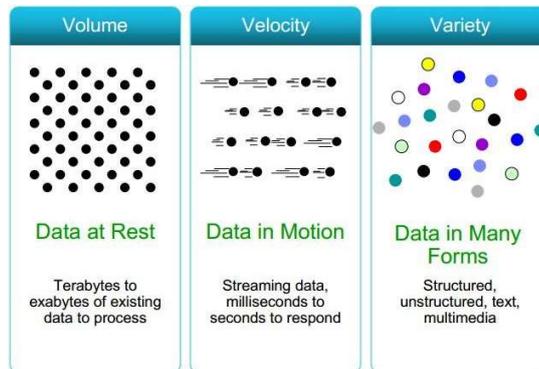
*Propostas por Doug Laney,  
Consultoria Gartner, 2001*

## Características de Big Data

Expansão em 3 eixos e em taxas crescentes

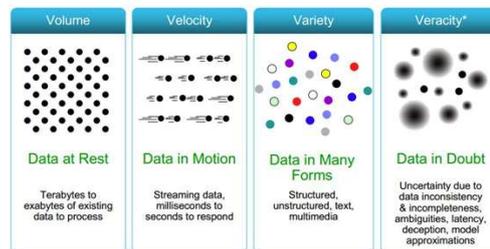


# Três Vs

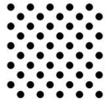
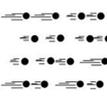
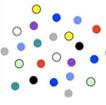


# Quatro Vs

Big Data Characteristics



# Cinco Vs

Volume	Velocity	Variety	Veracity	Value
				
<b>Data at Rest</b> Terabytes to Exabytes of existing data to process	<b>Data in Motion</b> Streaming data, requiring milliseconds to seconds to respond	<b>Data in Many Forms</b> Structured, unstructured, text, multimedia,...	<b>Data in Doubt</b> Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations	<b>Data into Money</b> Business models can be associated to the data

Adapted by a post of Michael Walker on 28/November 2017

# Vale o que é ouro

Tio Patinhas



## Vale o que é dado

- 1 tonelada de ouro vale US\$ 39,289,360.00
- Até hoje, 185.000 toneladas de ouro foram extraídos
  - Tudo somado: US\$ 7.4 trilhões (7400 000 000 000)
  - FAAMG + BAT: US\$ 5 trilhões (Junho de 2018)
    - Facebook, Amazon, Apple, Microsoft and Google
    - Baidu, Alibaba and Tencent



© André de Carvalho - ICMC/USP

29

## Vale o que é dado



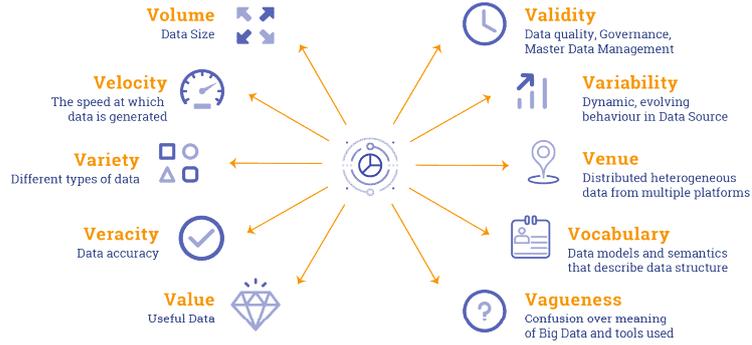
<https://www.nytimes.com/es/2017/05/12/amazon-apple-facebook-microsoft-y-google-nos-tienen-atrapados>

André P L F de Carvalho

30

# 10 Big Vs (2014)

## THE 10 Vs OF BIG DATA

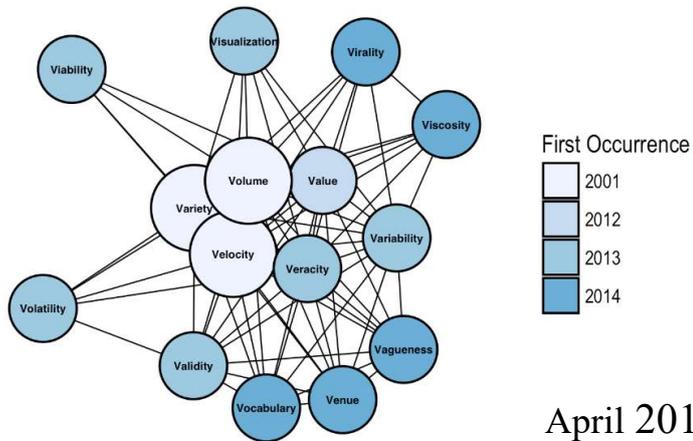


Source: xmonstack.com

© André de Carvalho - ICMC/USP

31

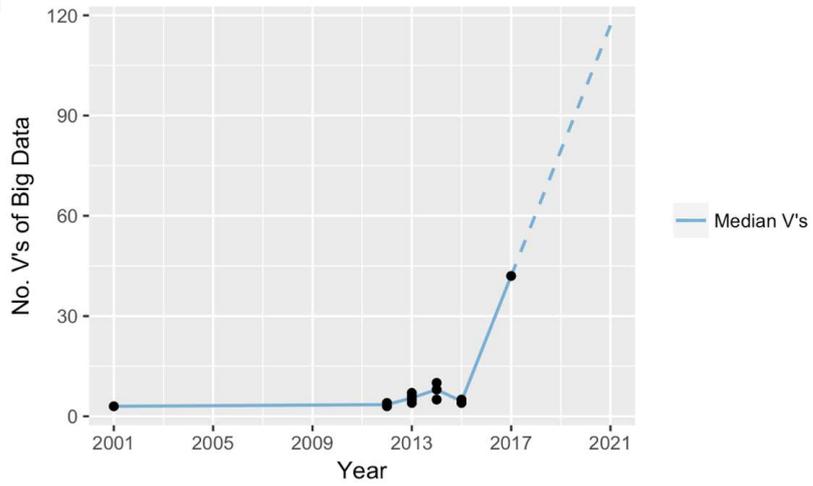
# 42 Big Vs



© André de Carvalho - ICMC/USP

32

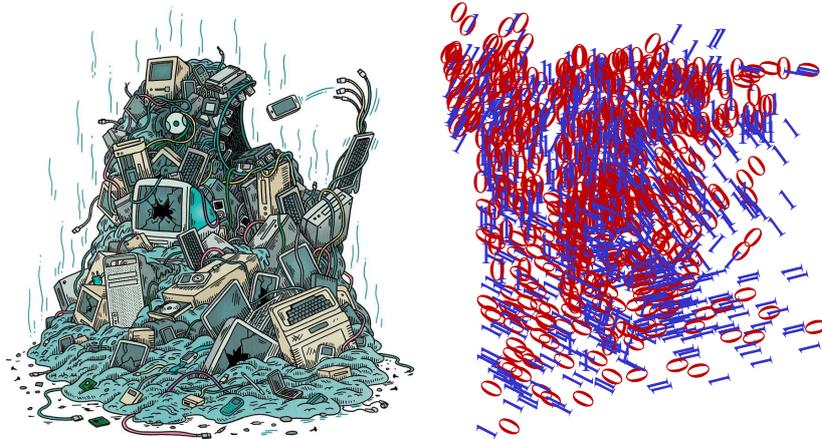
## Inflação de Vs



© André de Carvalho - ICMC/USP

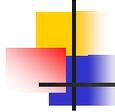
33

## Garbage Data



© André de Carvalho - ICMC/USP

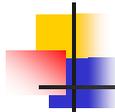
34



## Exploração de Big data

---

- Dados gerados geralmente contêm informações relevantes
  - Uma vez analisados, podem trazer vários benefícios
    - Sociais e econômicos
  - Crescente interesse na análise de dados



## Análise de dados

---

- Análise de dados por analistas
  - Faltam especialistas
  - Quando tem, o custo é elevado
  - Dificuldade de lidar de forma eficiente com grandes volumes de dados
- Necessário melhores ferramentas de análise
- Vários avanços nas últimas décadas
  - Ciência de Dados

## Ciência de Dados

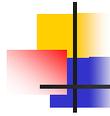
- Estuda princípios e técnicas para extrair conhecimento de um conjunto de dados
  - Novo, relevante e útil
- Pergunta chave da área:
  - Como extrair (de forma eficiente) conhecimentos em (grandes) conjuntos (fluxos) de dados



## Big Data x Ciência de Dados

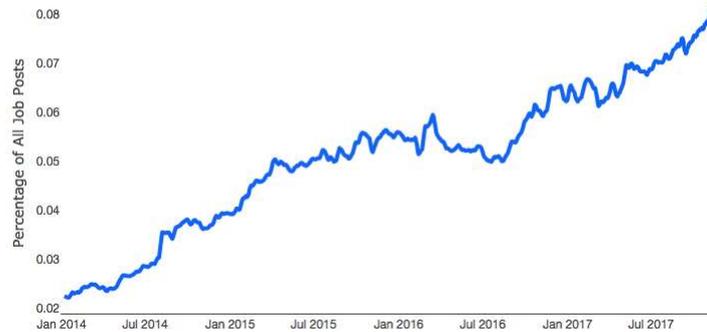
- Frequentemente usados como sinônimos
  - Big Data lida com tecnologias para coletar, gerenciar e processar (Big) dados
  - Ciência de Dados lida com criação de soluções para modelagem de dados
    - Capazes de extrair conhecimento de dados reais

### Processar x Descobrir



## Empregos em Ciência de dados

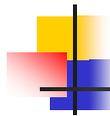
Data Science Jobs on Indeed Have Quadrupled over the last 4 years



<https://medium.com/indeed-data-science/transitioning-from-academia-to-industry-perspectives-from-indeeds-data-scientists-de890acd1bfc>

© André de Carvalho - ICMC/USP

39



## Quem esta contratando CD

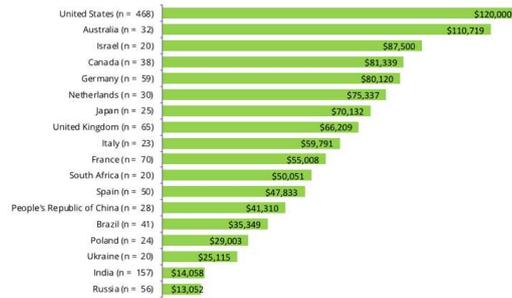
- Amazon
- Apple
- Electronic Arts
- Expedia
- Facebook
- Ford
- Google
- IBM
- Microsoft
- Booking.com
- Disney
- Ford
- Greepeace
- Hasbro
- Mattel
- Mercedes-Benz
- Nespresso
- Red Bull F1

© André de Carvalho - ICMC/USP

40

# Salários em Ciência de Dados

## Median Annual Compensation for Data Scientists and Machine Learning Engineers for Specific Countries



All values are in US Dollars. Conversion rates from mid 2017 (when the data were collected) were used for conversion. Data are from the Kaggle 2017 The State of Data Science and Machine Learning study. You can learn more about the study and download the data here: <https://www.kaggle.com/surveys/2017>. Countries are ranked by median annual salary. Only countries with ample sample size (n >= 20) are presented.



Copyright 2018 Business Over Broadway

© André de Carvalho - ICMC/USP

41

# Ciência de Dados

- Mineração de dados é o conceito mais semelhante a ciência de dados
    - Mas CD é muito mais que isso
  - Inclui
    - Compreensão do problema
    - Planejamento de experimentos
    - Pré-processamento
    - Modelagem
    - Avaliação e Validação
- } Mineração de dados

© André de Carvalho - ICMC/USP

42

## Aprendizado de Máquina (AM)

- Computadores podem aprender a realizar uma tarefa
  - Ao invés de serem explicitamente programados para isso
  - Aprendem a partir de um conjunto de dados a criar modelos (funções, hipóteses)
- Dezenas de milhares de algoritmos
  - Centenas de algoritmos novos a cada ano



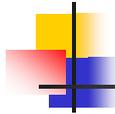
## AM vs computação tradicional

Escreve algoritmo detalhando como resolver um problema

Computação Tradicional

Aprende a resolver um problema observando dados coletados dele

Aprendizado de Máquina



## AM vs computação tradicional



Pessoas  
Programam!

Computação  
tradicional



Computadores  
programam!

Aprendizado  
de Máquina



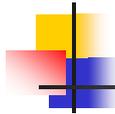
## AM vs computação tradicional

- Programação
  - Programa que funciona 90% das vezes é ruim
- Aprendizado de máquina
  - Modelo que acerta 90% das vezes pode ser o possível
    - E muito bom
    - Muitas vezes, é mais que suficiente para ser útil



## Aprendizado de Máquina

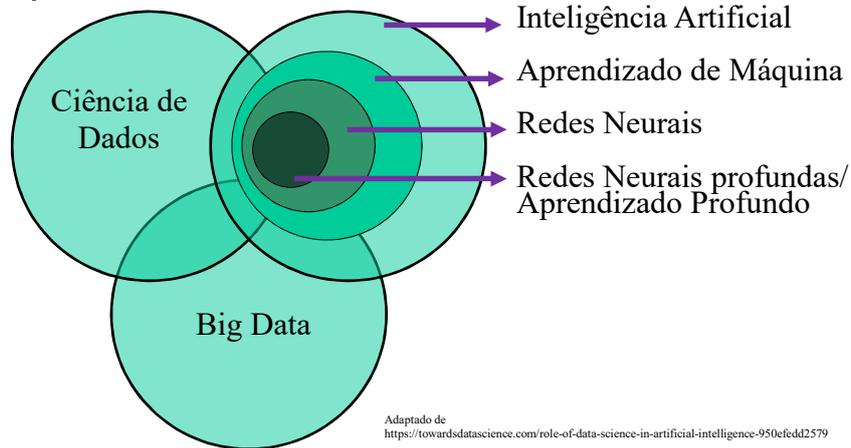
- Investiga algoritmos que possam aprender modelos a partir de dados
  - De forma automática, reduzindo (eliminando) intervenção humana
- Bem sucedido em vários problemas reais de modelagem
  - Descritivos
  - Preditivos



## Aplicações de AM

- AM está presente em várias atividades do nosso dia-a-dia
  - Recomendar que mensagens mostrar em aplicativos de redes sociais
  - Filtrar spams de seus *emails*
  - Decidir que resultados (e anúncios), e em que ordem, mostrar após uma busca na internet
  - Diagnóstico médico
  - Predizer aproveitamento de estudantes

## Como A, B e C se relacionam?



Adaptado de <https://towardsdatascience.com/role-of-data-science-in-artificial-intelligence-950efcdd2579>

Andre Ponce de Leon de Carvalho

49

## AM de ponta-a-ponta

### Utilizando Aprendizado de Máquina

Escolher/Modificar algoritmo de AM

Lidar com dados desbalanceados

Engenharia de atributos

Selecionar atributos



Ajustar hiperparâmetros

Lidar com valores ausentes

Verificar overfitting

Descobrir bugs

Adaptado de Rick Caruana, Research opportunities in AutoML. Microsoft Research



© André de Carvalho - ICMC/USP

50

# AM de ponta-a-ponta

## Aprendizado de Máquina na Prática

Escolher/Modificar algoritmo de AM

Lidar com dados desbalanceados

Engenharia de atributos

Selecionar atributos



Ajustar hiperparâmetros

Lidar com valores ausentes

Verificar overfitting

Descobrir bugs



Adaptado de Rick Caruana, Research opportunities in AutoML. Microsoft Research

© André de Carvalho - ICMC/USP

51

**Forbes** | Billionaires | Innovation | Leadership | Money | Consumer | Industry

**ZDNet** | VIDEOS | GC | WINDOWS 10 | CLOUD | AI | INNOVATION | SECURITY | MOBILE | NEWSLETTERS | ALL

MUST READ: These hackers broke into 10 telecoms companies to steal customers' phone records

### AutoML is democratizing and improving AI

Once a niche technology, Automated Machine learning (AutoML) is now a thing. Helping non-data scientists do simple AI, and helping trained data scientists do complex work ever-faster, AutoML technology is catching on, and may well put AI in the Enterprise fast lane.

## Why AutoML Is Set To Become The Future Of Artificial Intelligence

19,061 views | Apr 15, 2018, 02:05am

### 3) Say "Hello" To AutoML

One of the biggest trend that will dominate the AI industry in 2019 v automated machine learning (AutoML). With automated learning capabilities, developers will be able to tinker with machine learning and create new machine learning models that are ready to handle future AI

### The Rise of Automated Machine Learning

No matter what industry you're in, autoML can help you use machine learning successfully and leverage business insights hidden in places where only machine learning can reach.

By Abhi Yadav  
January 15, 2019

PRODUCTS - ARTICLES - SERIES - CATEGORIES - ABOUT

### AutoML - A Short Overview: Why AutoML Is Ready To Be The Future Of Artificial Intelligence

### AutoML - The Future Of Artificial Intelligence | CIS Coff...

Watch later | Share | cisin.com

### HEARTBEAT

MOBILE | MACHINE LEARNING | NEWSLETTER | COMMUNITY

### AutoML: The Next Wave of Machine Learning

Parul Pandey | Follow  
Apr 18 - 9 min read

### Love it or Hate It: Auto ML is Here to Stay

By David Sweeney - March 28, 2019

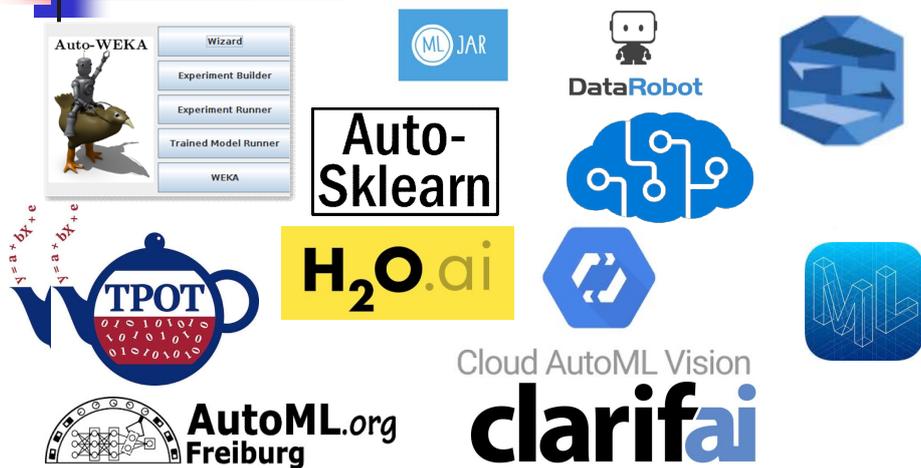
### AutoML and the Rise of Advanced Machine Learning Platforms

FEBRUARY 7 2019 - 1 COMMENT

Andre Ponce de Leon de Carvalho

52

## Ferramentas de AutoML



© André de Carvalho - ICMC/USP

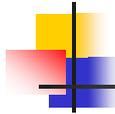
53

## CD Responsável

- Reprodutibilidade
  - Disponibilização e curadoria de dados e códigos
- Privacidade
  - Aplicando AM a 10 (300) likes, é possível conhece melhor sua personalidade que colega de trabalho (cônjuge)
  - *Fair Information Practices* (FIPs)
- Responsabilidade (*accountability*)
  - Alguém responde pela ferramenta de CD

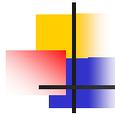
© André de Carvalho - ICMC/USP

54



## CD Responsável

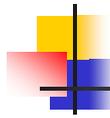
- Justiça
  - Tomada de decisão não deve embutir preconceito
- Transparência
  - General Data Protection Regulation (GDPR-UE)
    - Direito a informação
  - Brasil tem legislação semelhante



## Ciência de Dados para o Bem

- Movimento sem fins lucrativos
  - Levar benefícios sociais para as pessoas e comunidades
  - Alguns programas são adotados por empresas
- Como ele ocorre?
  - Reuniões
  - Eventos
  - Estágios acadêmicos
  - Redes sociais

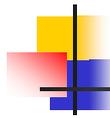




## Ciência de Dados no ICMC USP

- Encaixa perfeito com perfil do ICMC
- Pesquisa
  - Um dos 3 pilares do CEPID CeMEAI, sediado no ICMC
  - Centro de Excelência da Intel em Inteligência Artificial
- Extensão
  - Escolas de Ciência de Dados e de Big Data
  - MBA em Ciência de Dados

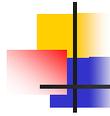
57



## Ciência de Dados no ICMC USP

- Graduação
  - Ênfase de graduação em Ciência de Dados
  - Curso de graduação em Ciência de Dados (trâmites legais iniciados)
- Pós-graduação
  - Mestrado Profissional

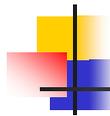
58



## Ciência de Dados no ICMC USP

- Mais de 20 docentes trabalhando em
  - Big Data
  - Ciência de Dados
  - Aprendizado de Máquina
- Laboratórios
  - Analytics, BioCom, LABIC, Estatística
- Formação de recursos humanos
  - Graduação e Pós-graduação

59

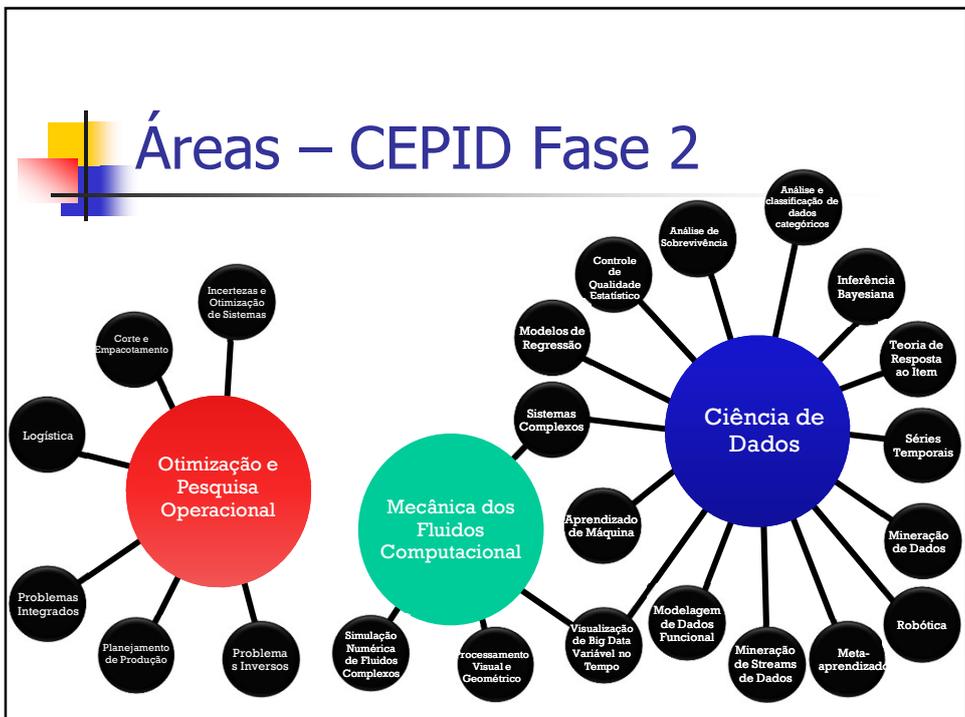
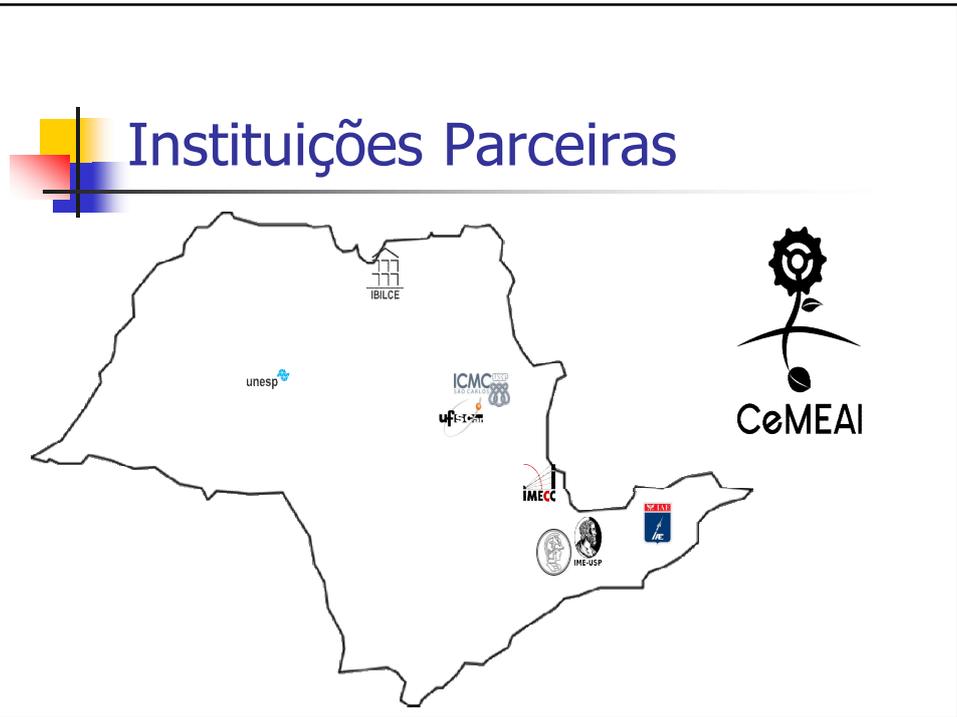


## CEPID CeMEAI



- Centro de Ciências Matemáticas Aplicadas à Indústria
  - ICMC-USP
- Início em 2013, duração de até 11 anos
  - Orçamento para primeiros 5 anos de 15 milhões de Reais da FAPESP
    - Mesma quantidade obtida de empresas
- Objetivo principal
  - Transferir conhecimento em computação, matemática e estatística para empresas e órgãos públicos

60



# Workshops com indústrias



2015  
7 problemas de 6 empresas  
139 participantes  
Resultados  
1 proposta de PIPE  
1 proposta proposta de cooperação  
2/3 problemas completamente resolvidos

2016  
6 problemas de 5 empresas  
137 participantes  
Resultados  
1 problema completamente resolvido  
1 problema com 4 soluções possíveis de patenteamento  
4 problemas com grande avanço para a solução

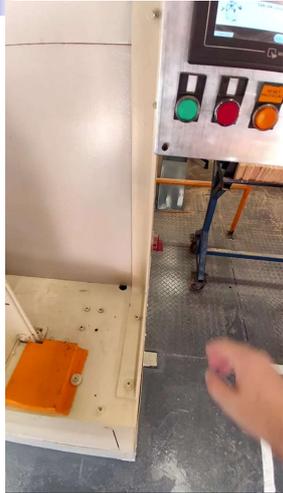
2017  
6 problemas de 6 empresas  
132 participantes  
Resultados  
2 problemas completamente resolvidos  
2 com grandes avanços

# Mestrado Profissional

6ª turma em Ciência de Dados  
Alunos de empresas:



# Classificação de madeiras



- Protótipo para Faber Castell
- 100 tábuas
  - 50% boas and 50% ruins
    - 14 tipos de defeitos
  - 150 placas por minute (27,7 metros por minute)
  - Acurácia preditiva: 97%
    - Operador humano: 93%

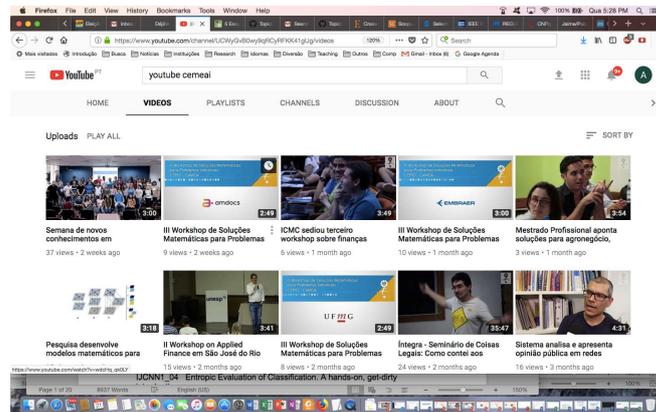
© André de Carvalho - ICMC/USP

# CEPID CeMEAI





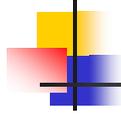
## Canal YouTube CeMEAI



67

## Conclusão

- Big data
- Ciência de dados
- Aprendizado de Máquina
- A, B e C no ICMC



# Perguntas

---

