

Introdução à Análise Exploratória de Dados

Noções Gerais sobre Data Science

João Pedro Albino

Departamento de Computação / Faculdade de Ciências

PPG-MiT / Faculdade de Artes, Arquitetura, Comunicação e Design



Oportunidade

Data, data everywhere

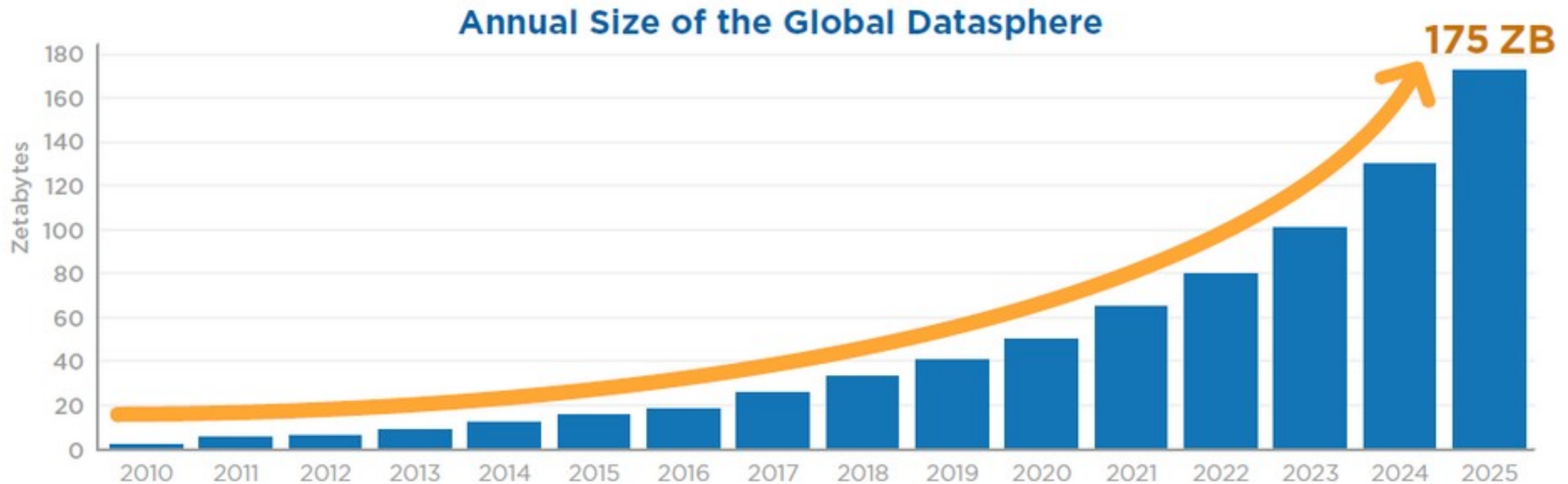
*Information has gone from scarce to superabundant. That brings huge *new* benefits, says Kenneth Cukier — but also big headaches.*

The Economist, Feb. 27th, 2010



Era do Big Data

Figure 1 - Annual Size of the Global Datasphere



O que é Big Data?

- Conjunto de dados brutos
- Coletados e armazenados, por vários meios
- Podem ser analisados computacionalmente
 - Revelar padrões
 - Tendências
 - Associações
 - Especialmente em relação ao comportamento humano e suas interações.



<https://www.integrity-ux.com.br/blog/2020/06/12/big-data-o-que-e-e-como-funciona/>

Os 5 V's do Big Data

- **Volume**

- É quantidade imensurável de dados que existe em todo o mundo hoje.

- **Velocidade**

- Com o avanço das tecnologias, a produção de dados é mais veloz e a tomada de decisão mais rápida torna-se cada vez mais importante.

- **Variedade**

- Com as diversas plataformas e meios de comunicação, as fontes de dados são mais variadas

- **Veracidade**

- Garantia de que os dados utilizados estão corretos e são válidos

- **Valor**

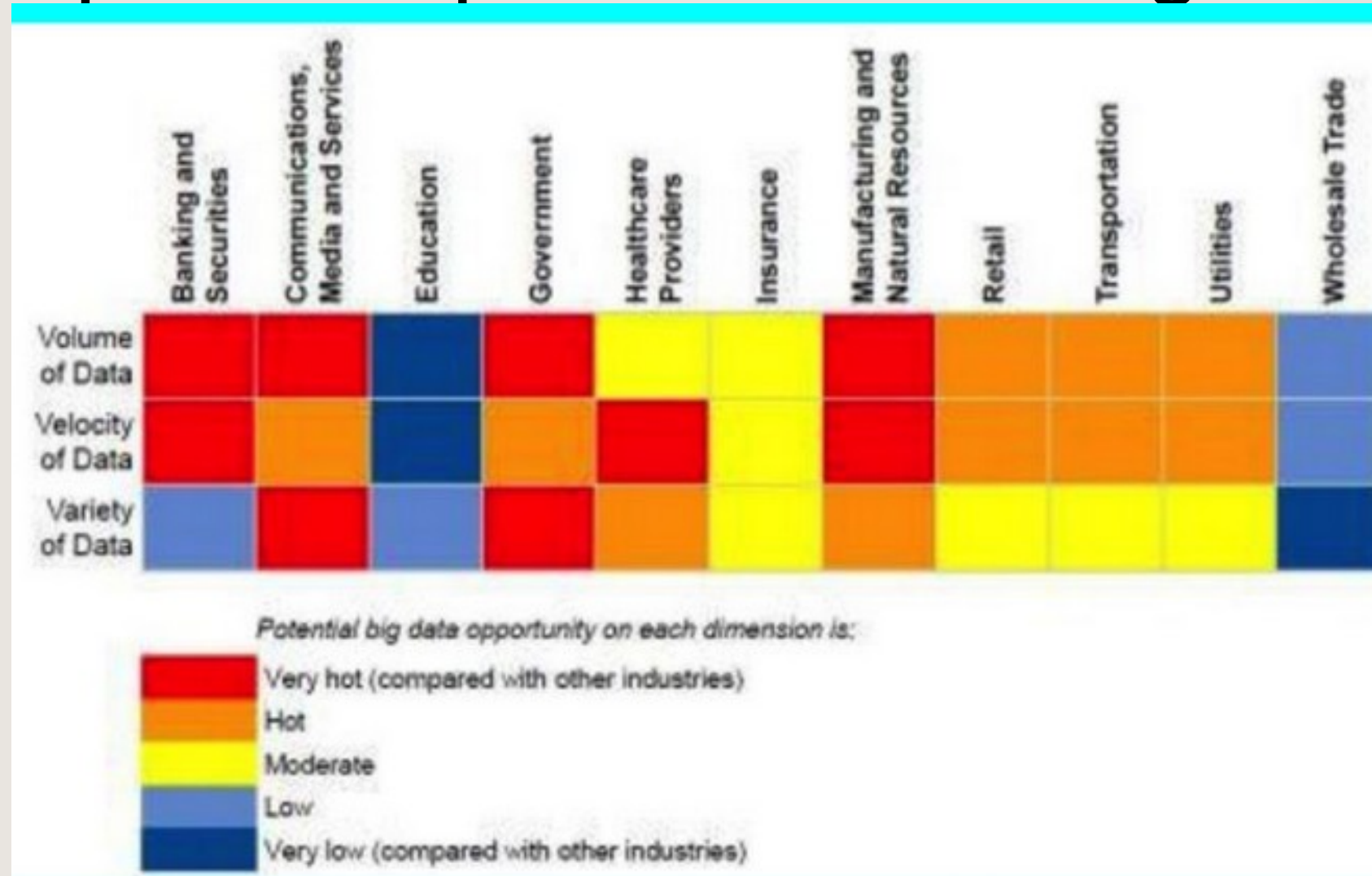
- Garantia de que os dados utilizados agreguem valor
- O valor que os dados geram para os usuários e para os negócios



Áreas de Aplicação do Big Data



Perspectivas para o uso de Big Data



Big Data: Principais Desafios



Desafios para o Big Data: Comunicações, mídia e entretenimento

- **Consumidores esperam rich media**
 - *Vídeo, áudio ou outros elementos que incentivam a interação e*
 - *Envolvimento **on-demand** em diferentes formatos e dispositivos*
 - *Coleta, análise e uso dos **insights** do consumidor*
 - *Aproveitar o conteúdo de mídia móvel e social*
 - *Compreensão dos padrões de uso de conteúdo de mídia em tempo real*
- **Aplicações de big data**
 - *Criar perfis detalhados de clientes*
 - *Criar conteúdo para diferentes públicos-alvo*
 - *Recomendar conteúdo sob demanda*
 - *Medir o desempenho do conteúdo*



Desafios para o Big Data: Educação

• **Usado de forma significativa no ensino superior**

- - medir eficácia do professor
 - garantir boa experiência para alunos / professores
 - desempenho do professor ajustado / medido
 - variáveis
 - número de estudantes, assunto, demografia estudantil, aspirações
- - estudantis, classificação comportamental

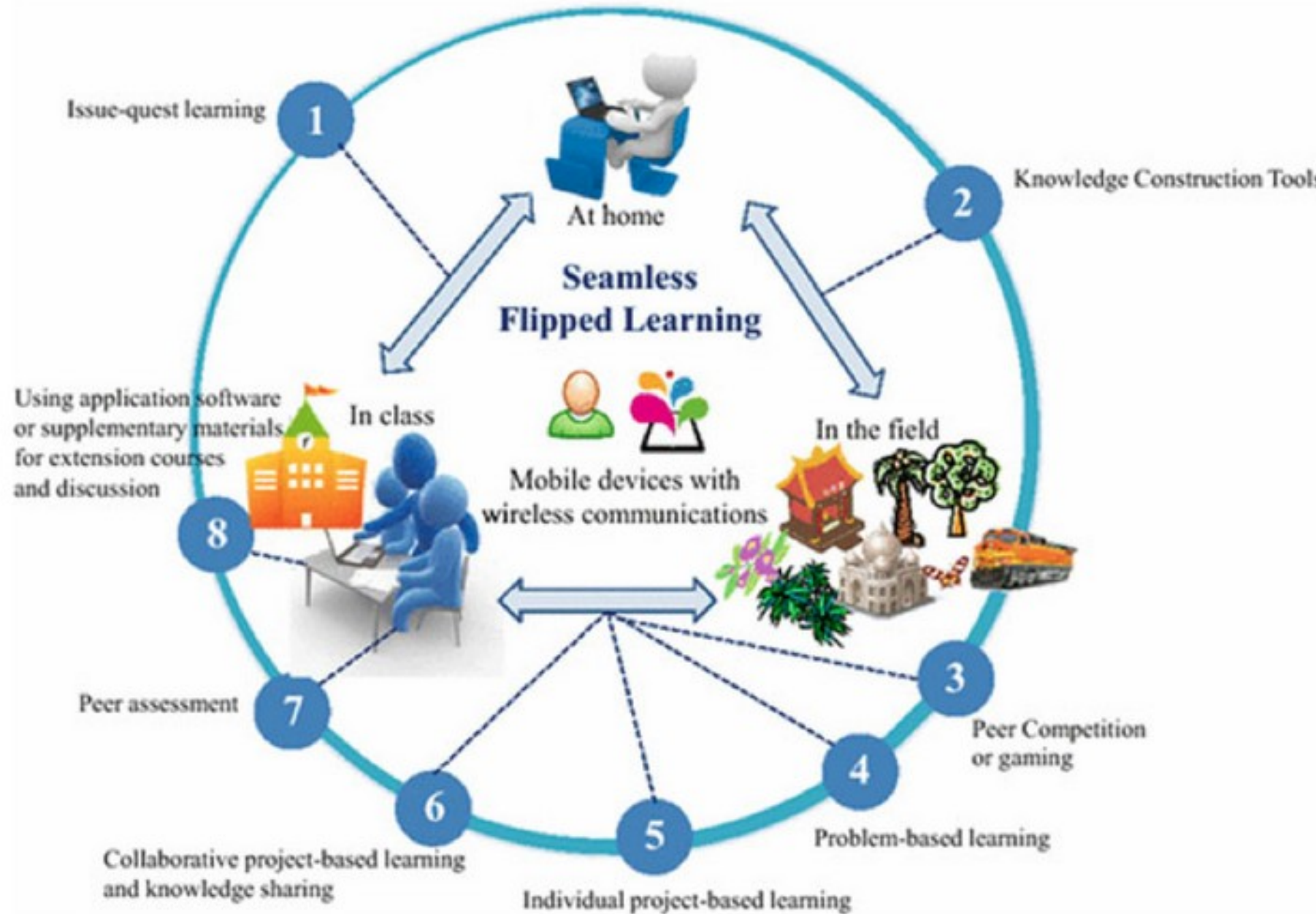
• **Departamento de Educação dos EUA**

- - grandes dados para análises
 - ajudar alunos que se distraem ao usar recursos on-line
 - padrões de clique usados para detectar aborrecimento



Desafios para o Big Data: Educação

- Seamless Learning



Desafios para o Big Data: Educação

- Learning **Analytics**



Desafios para o Big Data: no geral

- Qualidade e veracidade dos dados
- Segurança: manter os dados seguros ainda é um desafio
- Falta de profissionais com domínio na tecnologia e na área de aplicação
- Consciência do tipo de dado para análise
- Problema para armazenamento e gerenciamento dos dados
- Dificuldade em descobrir padrões e *insights*



Como trabalhar com Big Data



VOLUME



VARIETY



VELOCITY



VERACITY



VALUE

Ciência de Dados / Data Science

- **Definição**

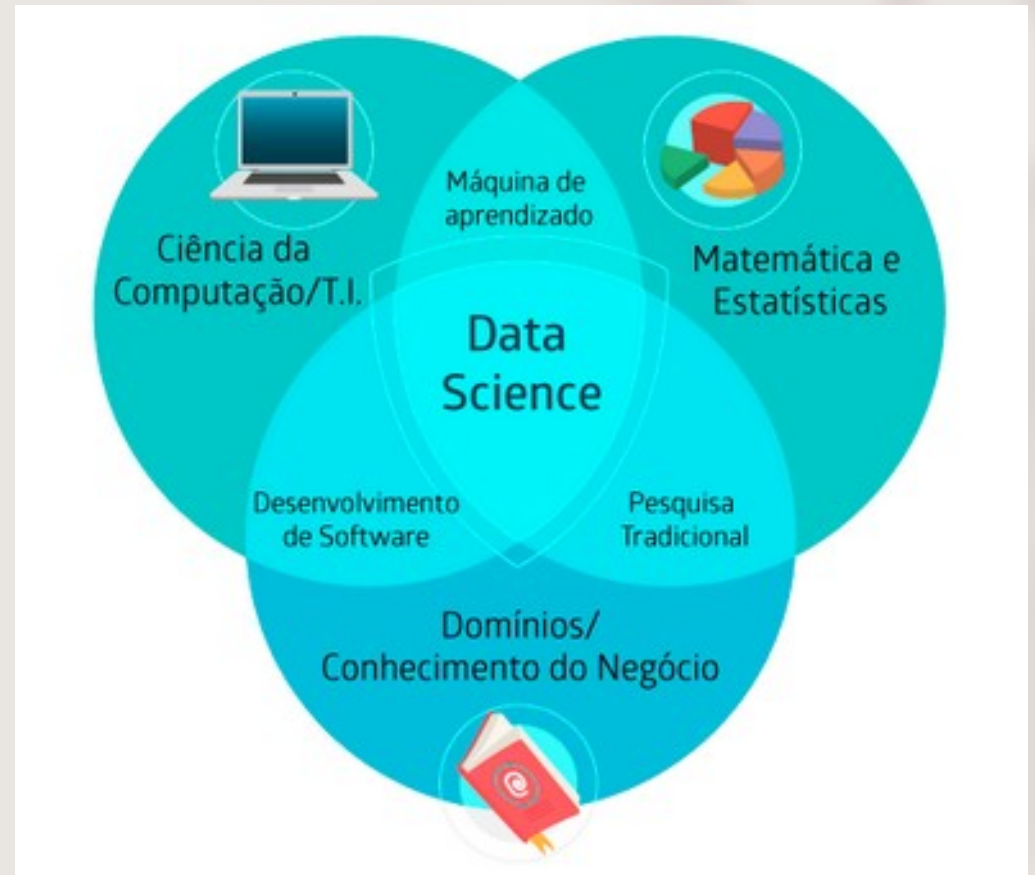
- *Termo utilizado para determinar uma combinação de várias ferramentas, algoritmos e princípios de aprendizado de máquina responsáveis pela descoberta de padrões e insights a partir de dados brutos.*
- Organizações enfrentam o desafio de lidar com enorme fluxo de informações geradas pela sociedade (big data)
- Saber utilizar a ciência de dados passou a ser um diferencial de negócios.
- É uma ferramenta valiosa para explorar e processar esses grandes volumes gerados por meio de diversas fontes.

O que é Data Science ?

Campo de estudo multidisciplinar que engloba dados, algoritmos e tecnologias capazes de extrair valor de dados estruturados ou não e resolver problemas analiticamente complexos.

É uma abordagem mais profunda, técnica e especializada sobre os elementos digitais.

Inclui o uso de modelos estatísticos e matemáticos, bem como de outras ferramentas para visualizar os dados.



Cientista de Dados: Profissão dos anos 20xx?

Review

Data | Data Scientist: The Sexiest Job of the 21st Century



Artwork: Tamar Co
on a page from a h

DATA

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

From the October 2012 Issue

Cientista de Dados: valores fundamentais



Um Cientista de Dados é alguém que sabe obter, tabular, explorar, modelar e interpretar dados, combinando e utilizando estatística e aprendizagem de máquina.



Cientistas de Dados não somente são adeptos a trabalhar com dados, mas apreciam esses dados como um produto de primeira classe.

Hillary Mason
Cientista de Dados
Accel



Habilidades de um(a) Cientista de Dados

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21st century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants

PROGRAMMING & DATABASE

- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing packages, e.g., R
- ☆ Databases- SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS

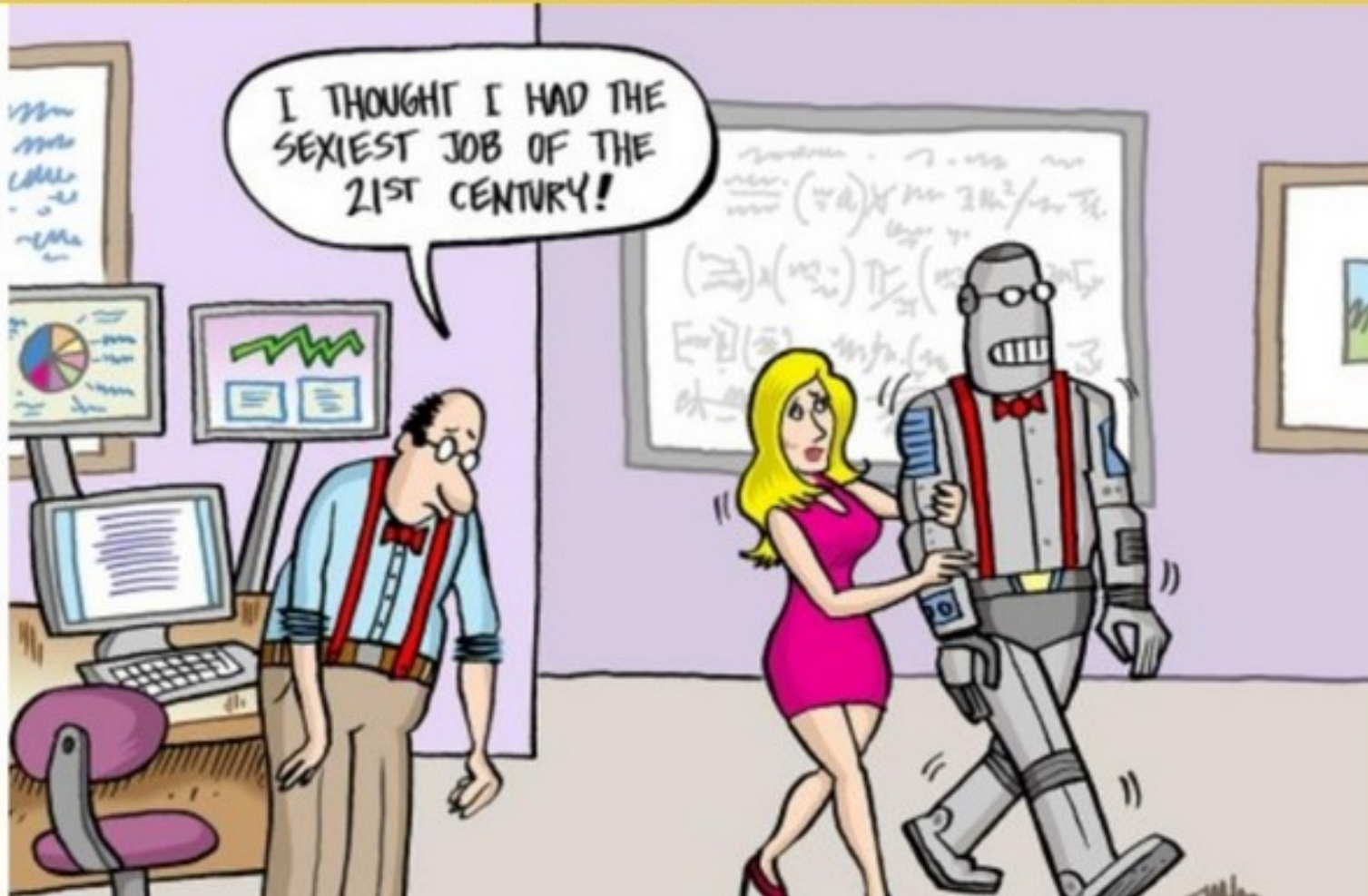
DOMAIN KNOWLEDGE & SOFT SKILLS

- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

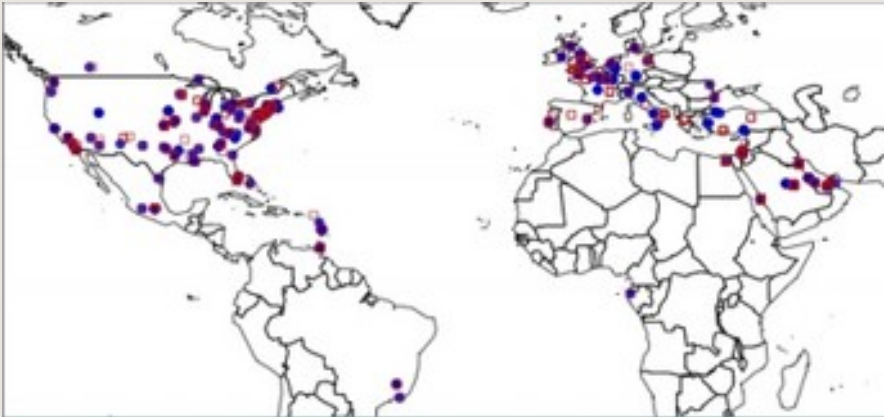
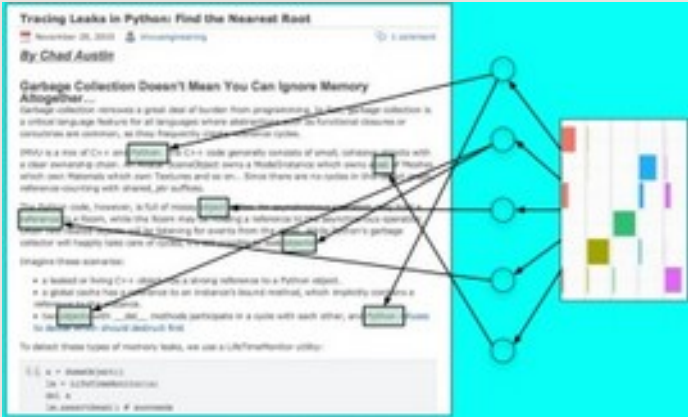
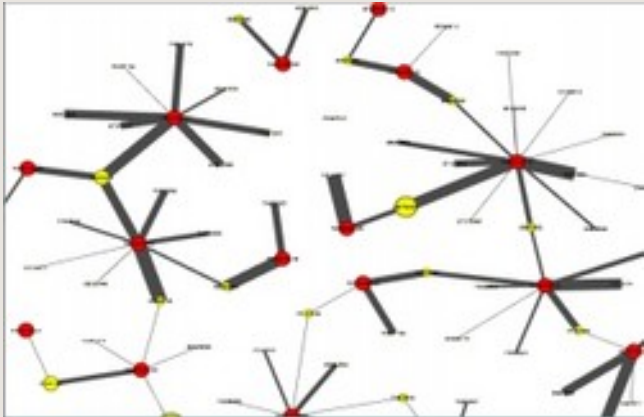
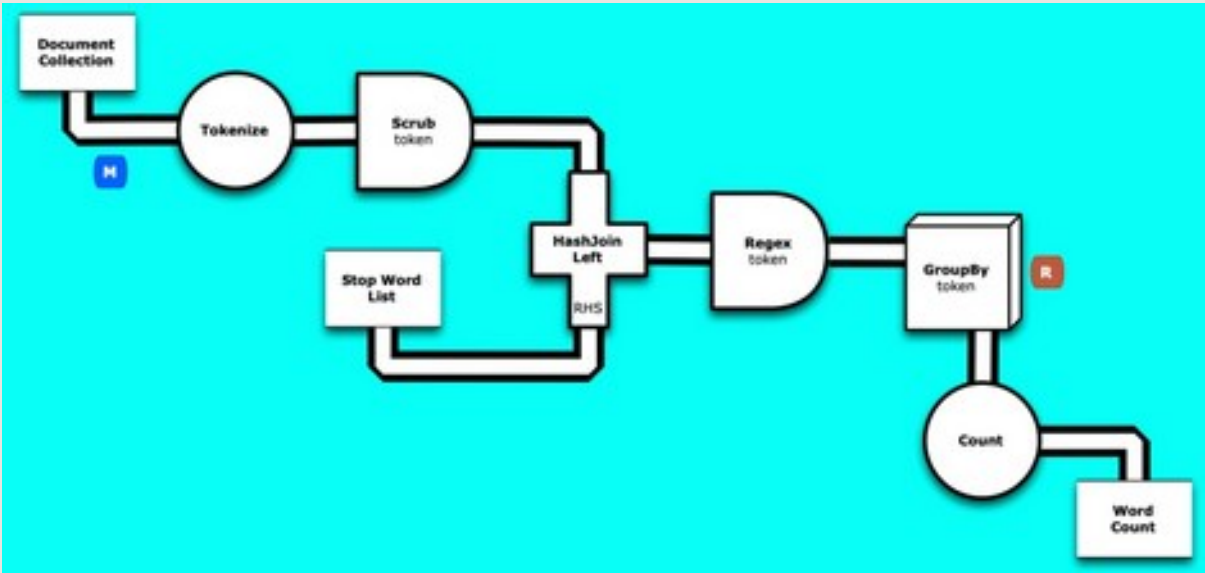
- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau





Cientista de Dados:
Realidade em
202x?

Ciência de Dados: bastidores



Pensamento Estatístico

um modo de raciocínio que inclui tanto o *raciocínio lógico quanto o analítico* :

- avaliar a **totalidade** de um problema, bem como suas **partes componentes**;

- procura avaliar os efeitos na mudança em uma ou mais variáveis!

esta abordagem tenta entender não apenas *problemas e soluções*, mas também os processos envolvidos e suas variações!

valioso no trabalho do Big Data quando combinado com experiência prática em comunicação

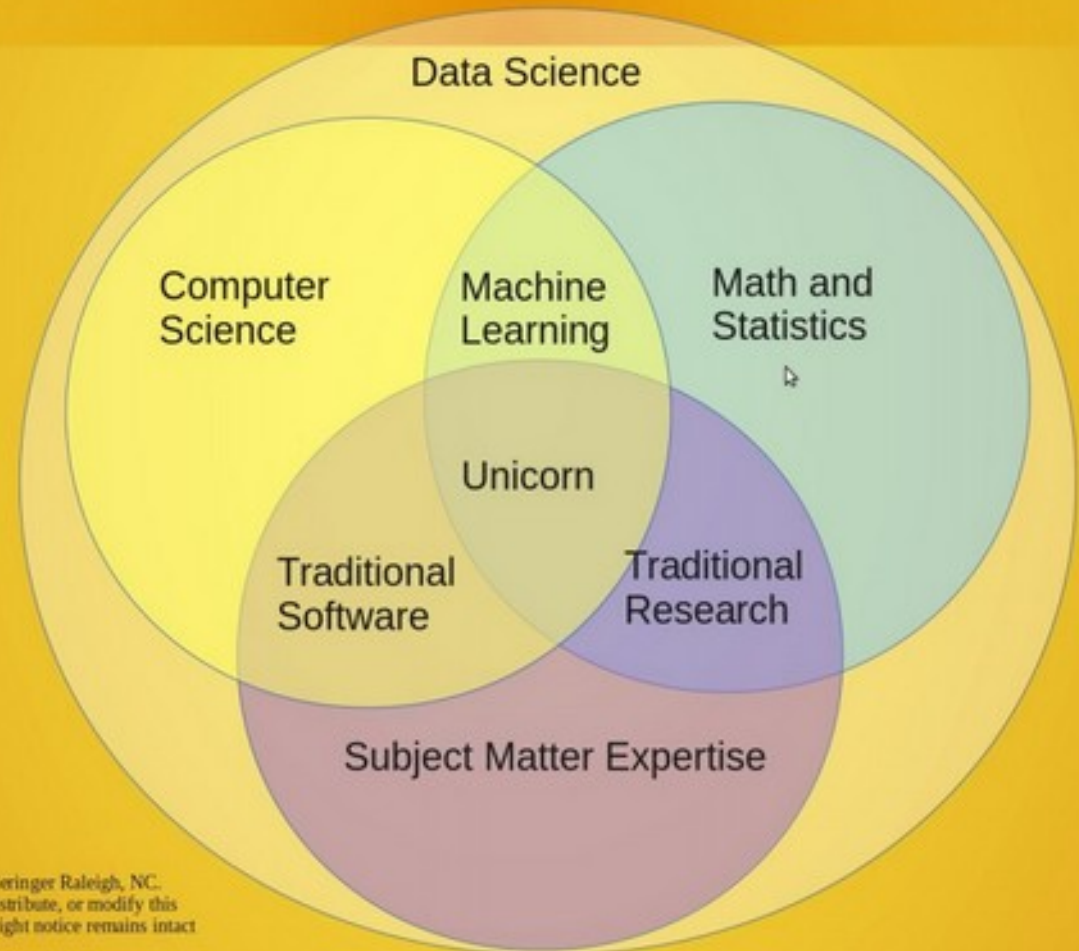
- p.ex.:aproximadamente 50% dos pares vêm do jornalismo ou de rádio e tv ...

- programadores normalmente não pensam assim!



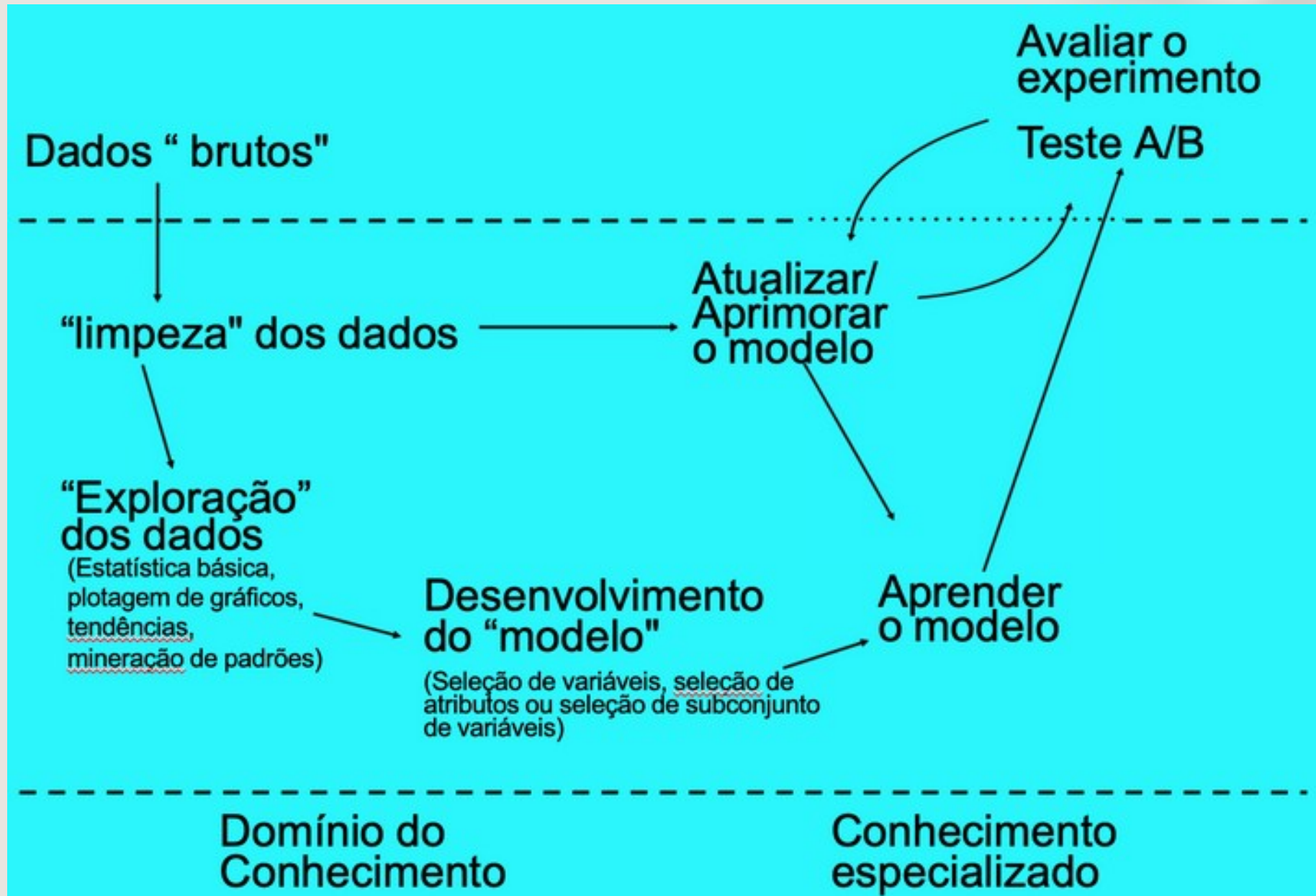
Disciplinas em Ciência de Dados

Data Science Venn Diagram v2.0

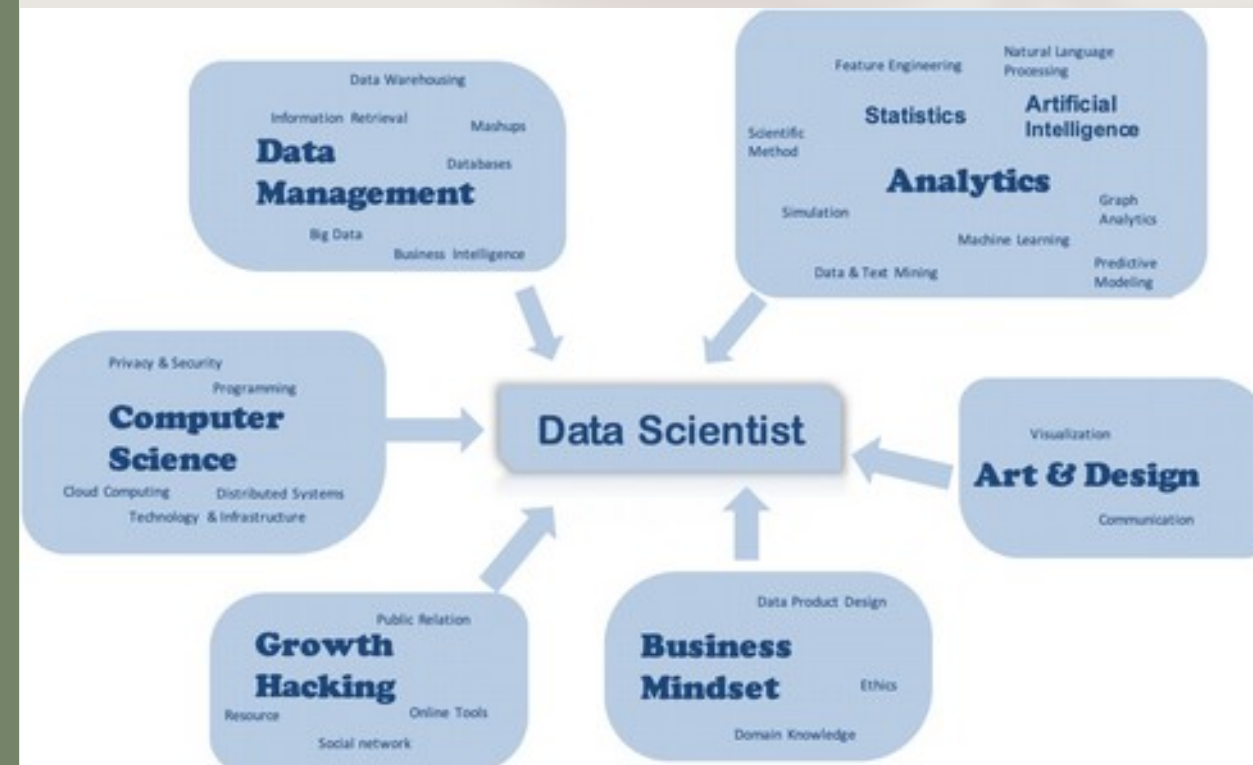


Copyright © 2014 by Steven Geringer Raleigh, NC.
Permission is granted to use, distribute, or modify this
image, provided that this copyright notice remains intact

Fluxo de Ciência de Dados



Quem é o cientista de dados?



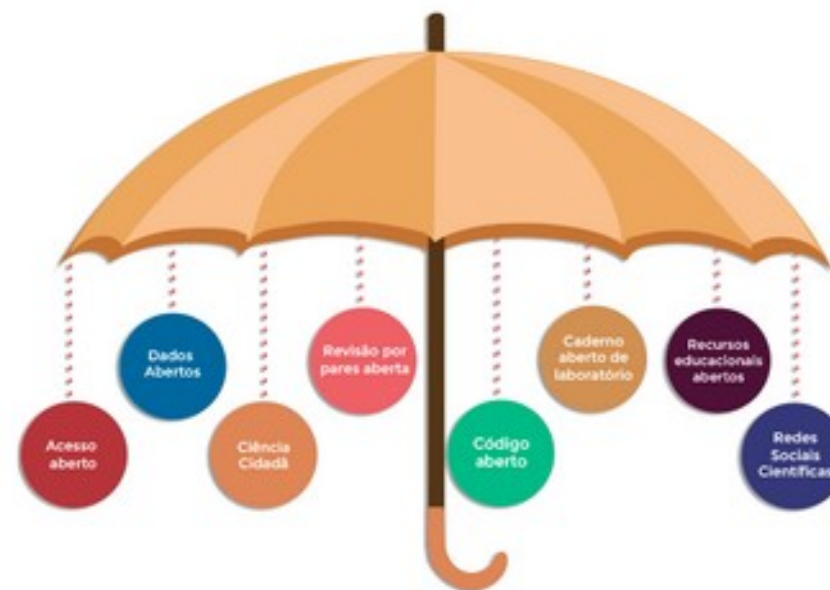
Algumas “pedras” no caminho

- **Data wrangling** ou **data preparation** – significa *preparação de dados*.
 - Conceito recente e diz respeito ao ato de **coletar**, **limpar**, **normalizar**, **combinar**, **estruturar** e **organizar** os dados que serão analisados.
- Aproximadamente 80% do custo dos projetos relacionados à análise de dados (**tempo + dinheiro**) são gastos na preparação destes dados (**data wrangling**) - principalmente nas questões de **limpeza** e garantia de **qualidade**.
- Infelizmente, os orçamentos relativos a dados tendem a entrar em frameworks que só podem ser utilizadas após a sua limpeza.



Ultrapassando obstáculos para obter “dados”

- “Criar” (gerar) os próprios dados, através de instrumentos (p.ex. questionários, surveys , etc.)
- Obter os dados através de APIs web, interfaces providas por bases de dados e por várias aplicações web modernas (incluindo Twitter, Facebook, dentre outras).
- Bases de dados "abertas", repositórios de dados, etc. (p.ex.: <https://www.cienciaedados.com/o-poder-do-open-data/>)
- Extrair as informações de arquivos (p.ex. PDF).
- Extrair informações de telas dos sites (scraping).



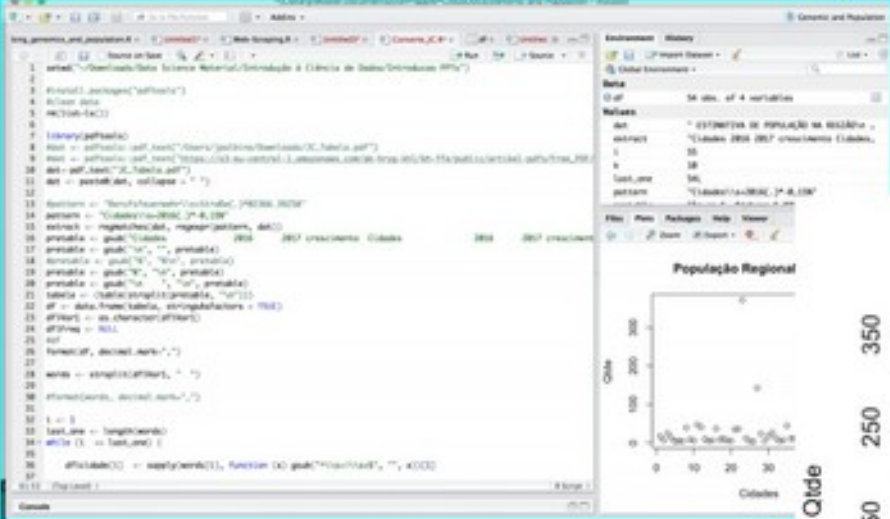
<https://portal.fiocruz.br/ciencia-aberta-na-fiocruz>

Exemplo



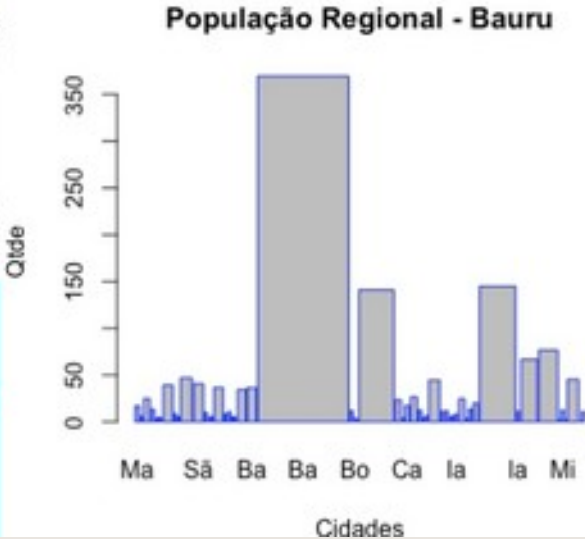
ESTIMATIVA DE POPULAÇÃO NA REGIÃO

Cidade	População
Baurópolis	371.690
...	...



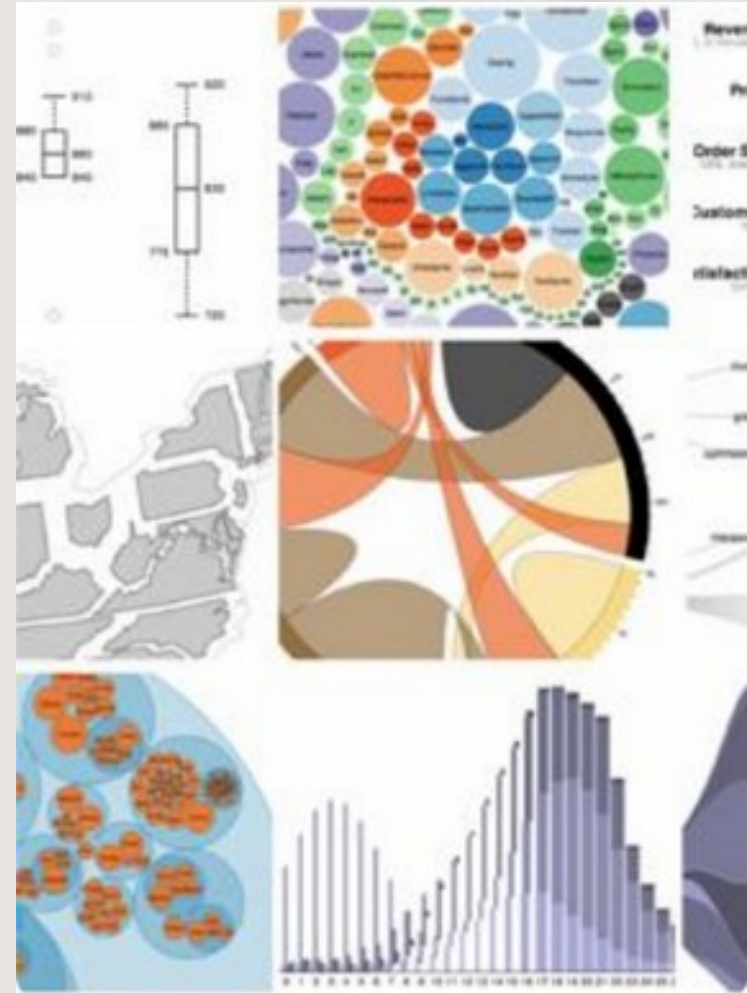
Excel

	A	B	C	D	E	F	G
1		cidade	pop2016	pop2017	indice		
2		Marcília	17.063	17.113	1,008		
3		Pardinho	6.195	6.239	1,009		
4		Praxi	24.762	24.979	1,009		
5		Praxi	11.114	11.135	1,002		
6		Pongá	3.494	3.48	0,99		
7		Pradópolis	1.076	1.126	1,046		
8		Pradópolis	4.147	4.215	1,019		
9		Pradópolis	26.113	26.126	1,001		
10		Pradópolis	8.844	9.042	1,022		
11		Santa	5.126	5.149	1,006		
12		Santa C. R.	46.893	47.148	1,005		
13		São Manuel	40.532	40.690	1,009		
14		Torrinha	8.89	8.934	1,004		
15		Ubatuba	4.711	4.734	1,005		
16		Agua S. M.	1.977	2.009	1,016		
17		Agua S. M.	36.704	36.88	1,005		
18		Aratiba	8.402	8.412	1,001		
19		Aratiba	11.078	11.107	1,002		
20		Aratiba	5.306	5.337	1,006		
21		Belém	5.006	5.188	1,036		
22		Baurópolis	34.528	34.602	1,002		
23		Baurópolis	36.526	36.510	0,999		
24		Baurópolis	369.346	371.690	1,006		
25		Baurópolis	11.026	11.04	1,001		
26		Baurópolis	4.675	4.717	1,009		
27		Baurópolis	1.548	1.577	1,019		
28		Baurópolis	141.052	141.546	1,004		

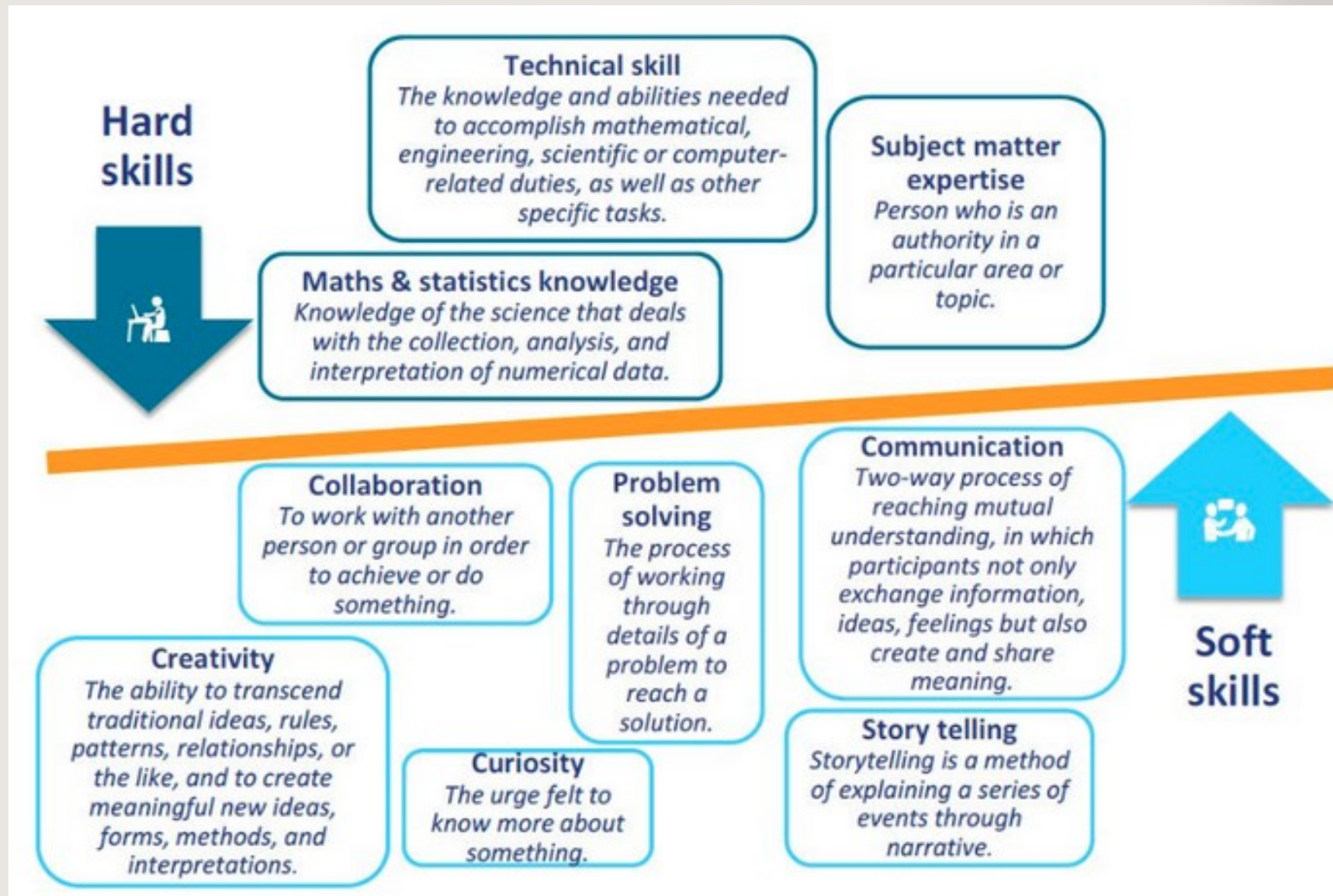


Habilidades mais importantes

- Aprender a utilizar ferramentas programáveis que preparem os dados
- **Aprender a gerar visualizações de dados convincentes**
- Aprender a estimar a confiança dos resultados relatados
- Aprender a automatizar o trabalho, tornando a análise repetível
- Outras importantes habilidades
 - Modelagem
 - Algoritmos
 - Etc.



Balancear “hard skills” e “soft skills”



Porque os cientistas de dados necessitam de ferramentas de visualização?



Porque os cientistas de dados necessitam de ferramentas de visualização?

Dar sentido aos dados ganhou grande importância neste século 21.

Programar é uma das formas de se manusear os dados disponíveis e torná-los utilizável.

Entretanto, nem todo mundo é “programador” e para tais pessoas os cientistas de dados podem usar ferramentas de visualização de dados para “contar histórias”.



Visualização de dados

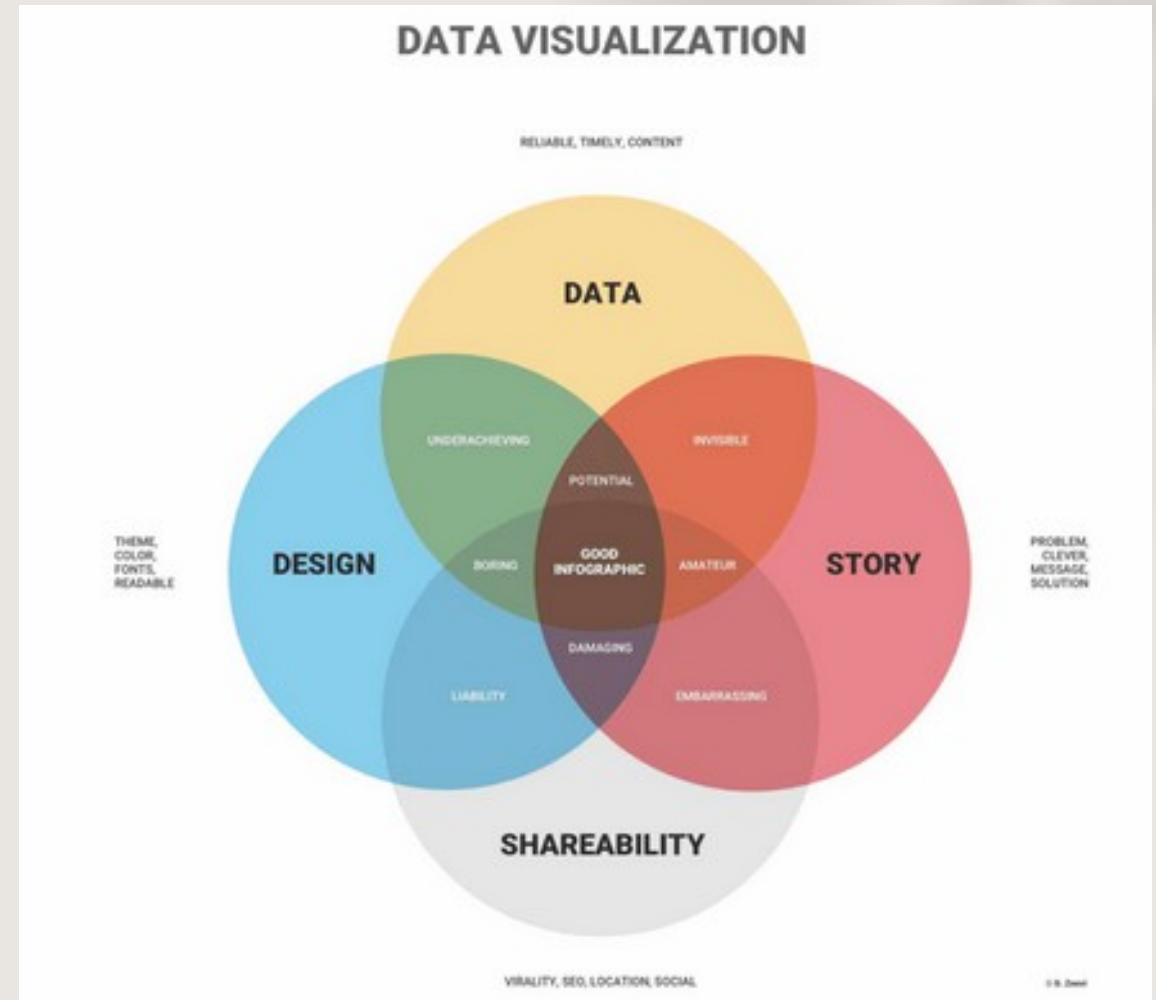
Visualização de dados é a apresentação de informações quantitativas em uma forma gráfica.

As visualizações de dados transformam conjuntos de dados (datasets) em visuais para o cérebro humano entender e processar.



Boas visualizações contam boas “estórias”

- Boas visualizações: **comunicação + ciência de dados + design**
- Utilizadas para descobrir fatos e tendências desconhecidos
- Entregam informações importantes sobre conteúdos complicados de *datasets* de maneiras significativas e intuitivas
- Eduard Tufte
 - "Ideias complexas comunicadas com clareza, precisão e eficiência"



Por que a visualização de dados (dataviz) é importante?

Melhor tomada de decisão

Organizações fazem melhores perguntas e tomam melhores decisões baseadas em dados.

Ênfase em monitoramento e desempenho: painéis de dados (dashboards) e KPIs.

Storytelling Significativo

Data Storytelling: técnicas que orientam a apresentação de informações e insights de dados

Gráficos de informações (infográficos): essencial para a grande mídia atual no jornalismo de dados.

Empresas: Cientistas de dados apresentam análises dos dados e os resultados utilizando data storytelling

Marketing: combinação de dados de qualidade e narrativa emocional

Alfabetização de dados

Ser capaz de compreender e ler as visualizações de dados tornou-se um requisito necessário para o século XXI.



: <https://www.datatelling.eu/what-is-data-storytelling/>



Inovação

Um(a) grande cientista de dados deve ser inovador(a) e criativo(a) com as habilidades que possui.

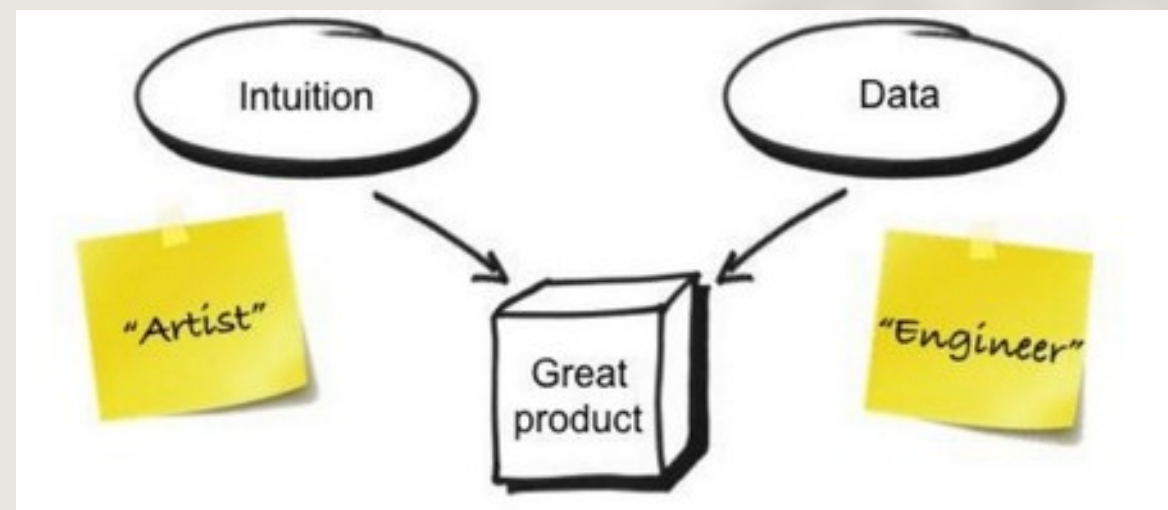
A criatividade de um(a) cientista de dados o(a) ajuda a determinar onde os dados podem agregar valor e trazer resultados lucrativos para a sua “aplicação”.

Intuição sobre os dados

Um bom cientista de dados não é aquele que apenas adiciona todas as possíveis características em um modelo de aprendizagem de máquina e analisa os resultados.

A coisa mais importante que um cientista de dados deve fazer antes de alimentar um modelo de aprendizagem de máquina é pensar se aquele modelo tem sentido!

Um cientista de dados deve ter **intuição sobre os dados**!





If you torture
the data long
enough, it will
confess to anything.

Ronald Coase

<https://i.redd.it/ikdymainj6p21.jpg>

Introdução à Análise Exploratória de Dados

Noções Gerais de Data Science

João Pedro Albino

Departamento de Computação / Faculdade de Ciências

PPG-MiT / Faculdade de Artes, Arquitetura, Comunicação e Design

