

Introdução à Ciência de Dados como Prática na Pesquisa Acadêmica

Jadson Pessoa

Professor do DECON | Membro GAPE



Plano de Trabalho

Objetivos

Workflow em Ciência de Dados

Linguagens de Programação em Ciência de Dados

R, RStudio e Repositórios

Atualizar

Coleta e Tratamento

Visualização

Modelagem

Comunicação

O mundo *tidyverse*

Objetivos

- Apresentar qual o entendimento que temos hoje do termo “**Ciência de Dados**” ou, como ficou mais conhecido, **Data Science**, em inglês.

O que é Ciência de Dados?

Data science combines the scientific method, math and statistics, specialized programming, advanced analytics, AI, and even storytelling to uncover and explain the business insights buried in data [IBM](#).

Data science combines multiple fields, including statistics, scientific methods, artificial intelligence (AI), and data analysis, to extract value from data [Oracle](#).

Data science [...] allow you to turn raw data into understanding, insight, and knowled. ([Wickham and Grolemund \[2017\]](#))

O que é Ciência de Dados?

Data science combines the scientific method, math and statistics, specialized programming, advanced analytics, AI, and even storytelling to uncover and explain the business insights buried in data [IBM](#).

Data science combines multiple fields, including statistics, scientific methods, artificial intelligence (AI), and data analysis, to extract value from data [Oracle](#).

Data science [...] allow you to turn raw data into understanding, insight, and knowled. ([Wickham and Grolemund \[2017\]](#))

O que é Ciência de Dados?

Data science combines the scientific method, math and statistics, specialized programming, advanced analytics, AI, and even storytelling to uncover and explain the business insights buried in data [IBM](#).

Data science combines multiple fields, including statistics, scientific methods, artificial intelligence (AI), and data analysis, to extract value from data [Oracle](#).

Data science [...] allow you to turn raw data into understanding, insight, and knowled. ([Wickham and Grolemund \[2017\]](#))

O que é Ciência de Dados?

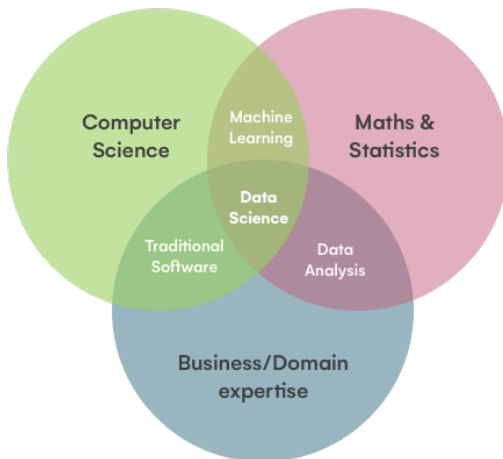


Figure 1: Big Data e Data Science

Big Data

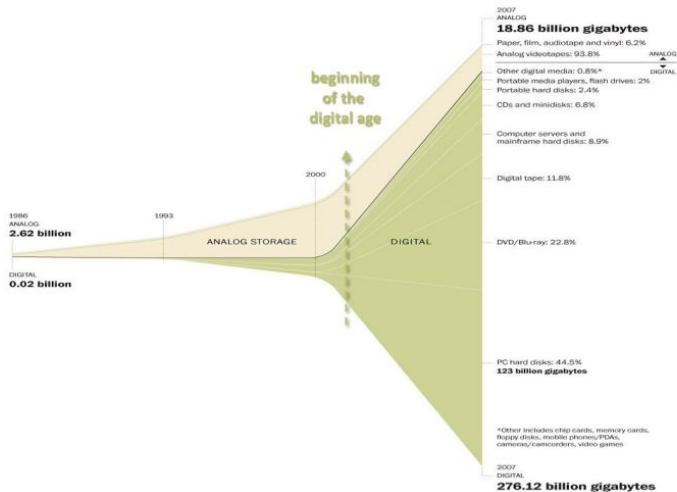


Figure 2: Big Data (Hilbert and López [2011])

Workflow em Ciência de Dados

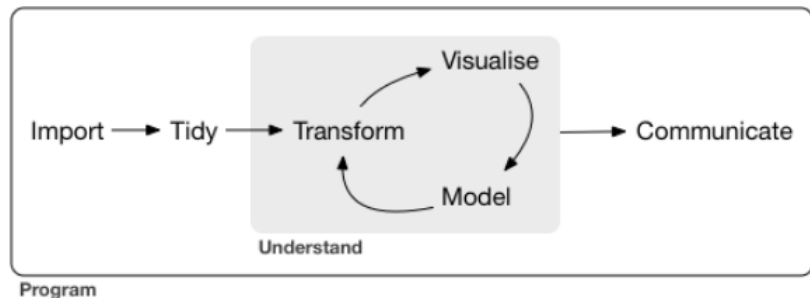


Figure 3: O processo de compreensão dos dados ([Wickham and Grolemund \[2017\]](#))

Linguagens de Programação em Ciência de Dados

- Em todas as etapas temos a utilização linguagem, ou melhor *linguagens*, de programação;

Atualizar

- Ciência de Dados (*Data Science*);
- Conjunto de técnicas de coleta, tratamento, análise e apresentação de dados;
- Sistematização de dados em busca informações úteis.
- **melhorar apresentação inicial**

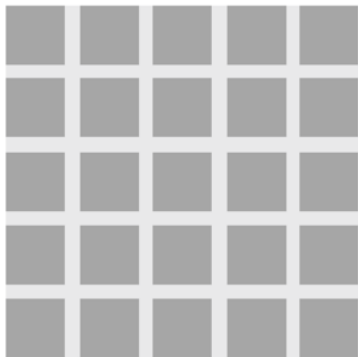
Introdução

- Sistematização de projeto em DS:
 1. Coleta;
 2. Tratamento;
 3. Análise;
 4. Visualização e apresentação.

Introdução

Coleta e Tratamento

Structured data



Unstructured data



Figure 4: Dados Estruturados e Não estruturados

Coleta e Tratamento

A **coleta** dos dados podem em diferentes formatos:

- Excel;
- XML;
- JSON;
- txt;
- HTML;
- MySQL;
- Formatos proprietários (Stata, Minitab, SPSS, SAS, etc).

Coleta e Tratamento

O dados precisam ser **tratamentados**:

- Limpeza de dados;
- Tratamento de *missing values* ou NA;
- Construção de números índices;
- Deflacionar valores correntes;
- Obtenção de taxas de crescimento;
- Tratando tendências;
- Dessazonalização;
- Subconjuntos (*subsetting*);
- Classificação;
- Utilização de *lags*.

Visualização

Uma vez que seus dados estejam arrumados, podemos passar para a parte de **exploração dos dados**. A exploração de dados é a arte de analisar seus dados, gerando hipóteses rapidamente, testando-os rapidamente, repetindo-os várias vezes. O objetivo da exploração de dados é gerar muitos leads promissores que você poderá explorar mais tarde com mais profundidade. Em geral, faz-se exploração de dados por meio da *visualização* desses dados.

Um bom processo de visualização de dados permite que possamos nos concentrar naquilo que realmente importa, deixando de lado relações não tão importantes.

Modelagem

Uma vez que tenhamos conseguido propor uma *hipótese de trabalho* através da etapa de exploração/visualização de dados, o próximo passo é propor um **modelo** entre as variáveis do nosso conjunto de dados. O objetivo da modelagem é capturar a essência de um conjunto de dados.

Comunicação

A última etapa do processo de *data science* é comunicar os resultados para clientes, gestores ou demais interessados. É uma fase absolutamente crítica do projeto. Isto porque, ao menos que você consiga se comunicar com a sua audiência, de nada valeu todo o trabalho realizado nas etapas anteriores.

O mundo *tidyverse*

De modo a fazer cada uma dessas etapas dentro do R, nós vamos utilizar a família de pacotes *tidyverse*. Assim, antes de qualquer coisa, certifique-se que você tenha o tenha instalado na atual versão do R.

```
install.packages('tidyverse')  
require(tidyverse)
```

A família *tidyverse*

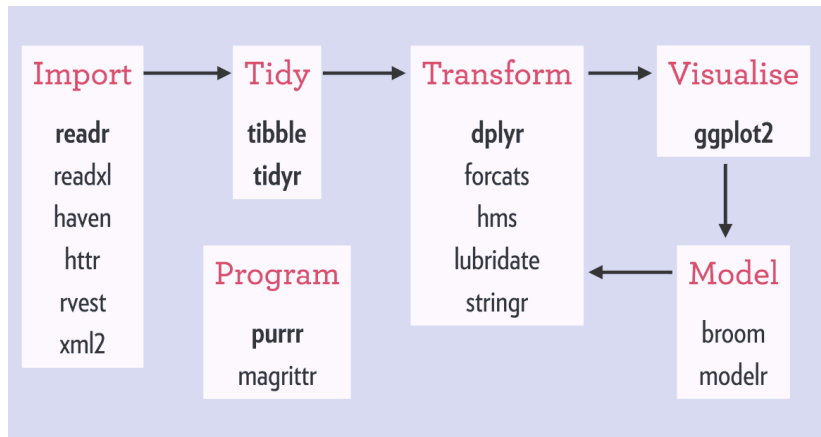


Figure 5: Os pacotes tidyverse (Wickham and Grolemund [2017])

Referências I

Martin Hilbert and Priscila López. The world's technological capacity to store, communicate, and compute information. *science*, 332(6025):60–65, 2011.

H. Wickham and G. Grolemund. *R for Data Science*. O'Reilly Media, 2017.