

# INF 1514

# Introdução à Análise de Dados

Material 1

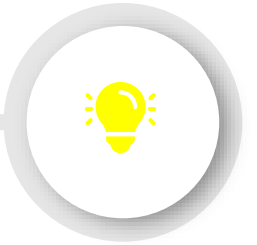


Este curso foi idealizado para **ensinar programação** para **análises básicas** e **visualizações de dados**.

INF 1514

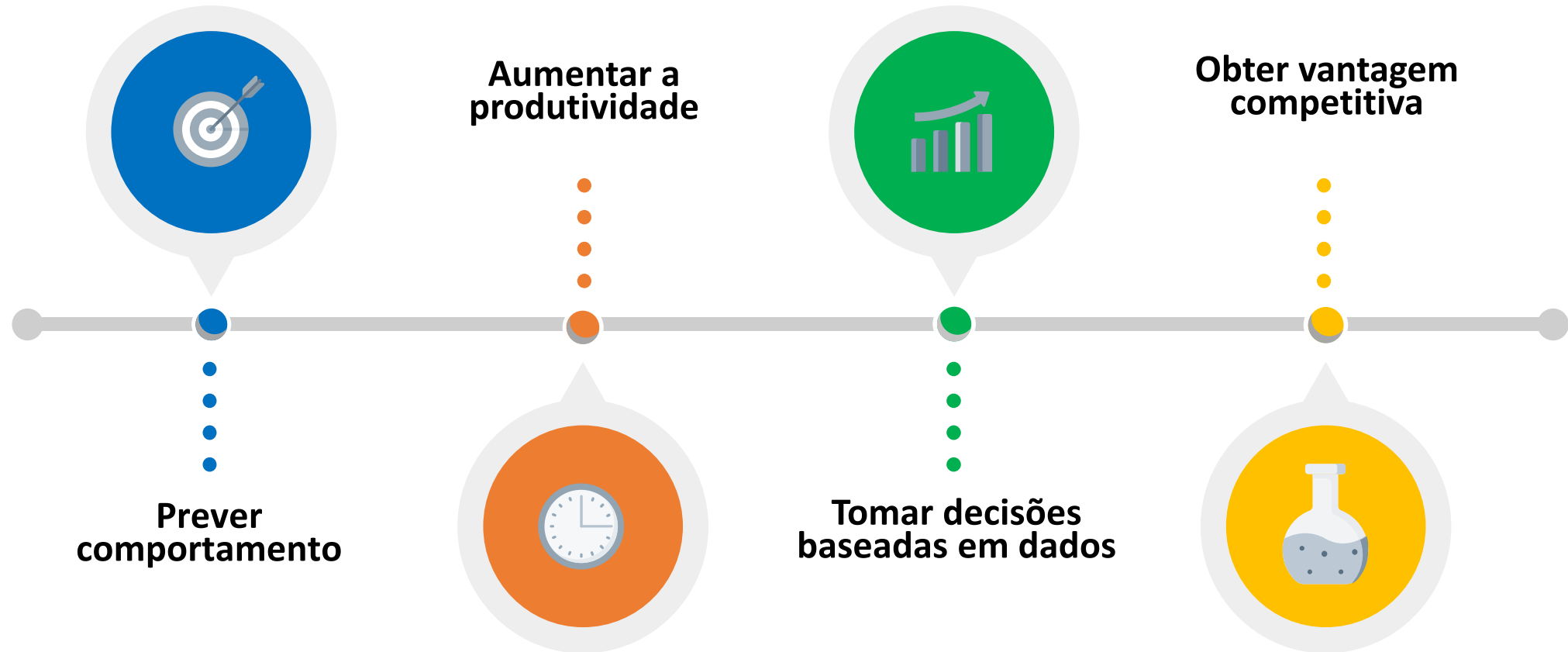


A **análise de dados** é um **processo** de inspeção, limpeza, transformação e modelagem de dados com o objetivo de **descobrir informações úteis, informar conclusões e apoiar a tomada de decisões.**



- Dados são simples observações sobre os estados do mundo.
- A análise de dados apresenta múltiplas abordagens, abrangendo diversas técnicas sob uma variedade de nomes, sendo usada em diferentes domínios de negócio, ciências e ciências sociais.
- Um exemplo simples de análise de dados pode ser visto sempre que tomamos uma decisão em nosso dia a dia, avaliando o que aconteceu no passado ou o que poderá acontecer no futuro se tomarmos uma determinada decisão.

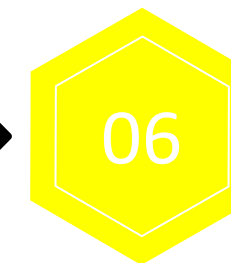
# Por que a Análise de Dados é importante?



# Fases do processo de Análise de Dados

## Definição dos requisitos de dados

Por que está fazendo essa análise? Quais dados irá utilizar?



## Limpeza dos dados

Os dados podem conter registros duplicados, espaços em branco ou erros.

## Interpretação dos dados

Interpretar os resultados e propor ações com base nas descobertas obtidas.

## Coleta dos dados

Coletar os dados com base nos requisitos, processá-los e organizá-los para análise.

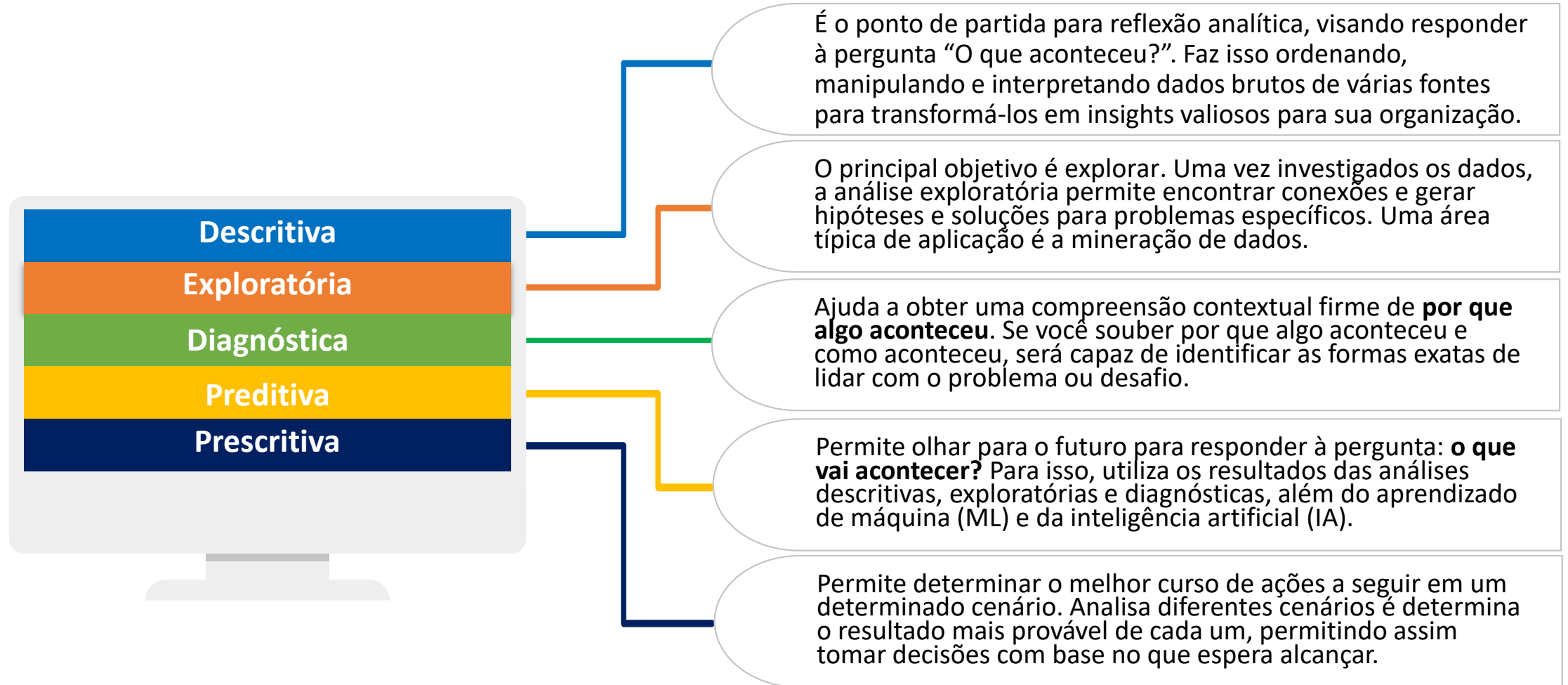
## Análise dos dados

Entender e interpretar os dados e tirar conclusões com base nos requisitos.

## Visualização dos dados

Apresentar graficamente as informações de forma que possam ser entendidas.

# Análise de Dados



# Métodos para Análise de Dados

1

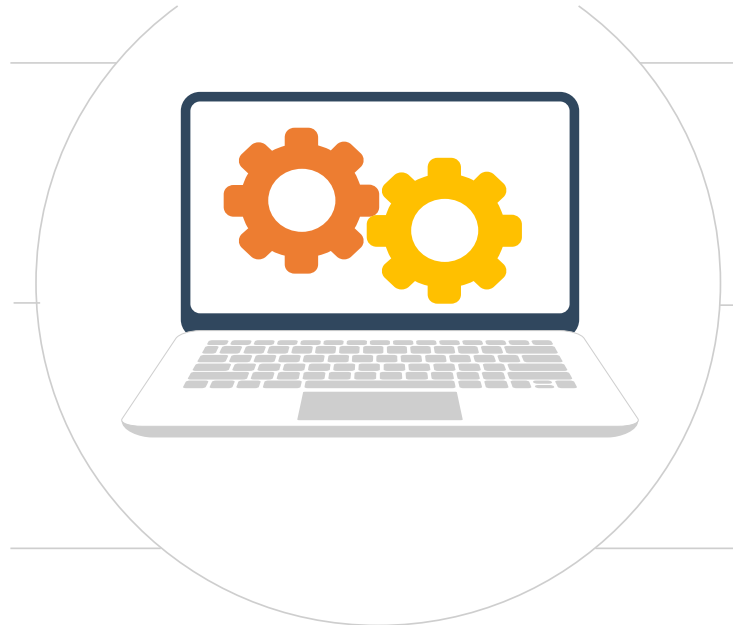
**Análise de Cluster**

2

**Análise de Regressão**

3

**Redes Neurais**



**Data e Text Mining**

4

**Árvores de Decisão**

5

**Análise de Séries Temporais**

6

# Ferramentas para Análise de Dados



Ajudam os usuários a processar, manipular e visualizar dados, analisar as relações e correlações entre conjuntos de dados e também a identificar padrões e tendências para interpretação.



# Saber programar é importante

1

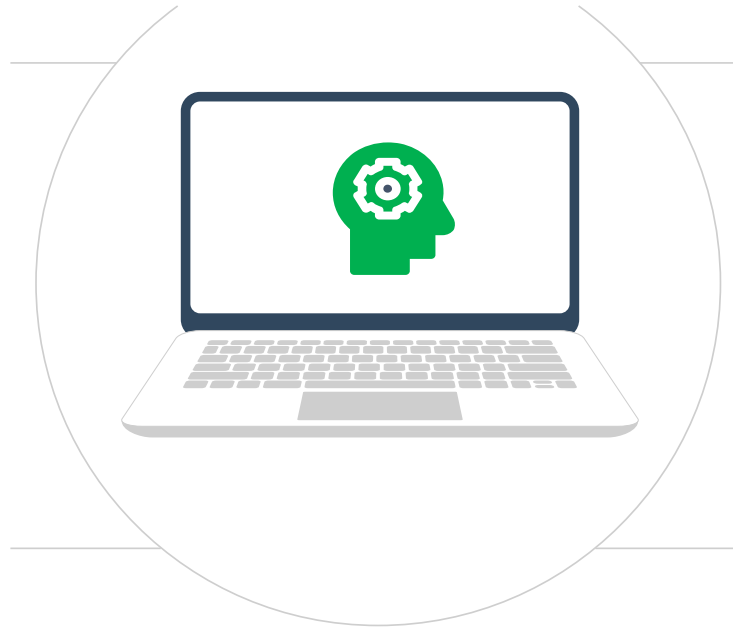
## Raciocínio lógico

Ajuda no desenvolvimento do raciocínio crítico, analítico e lógico.

2

## Habilidade de abstração

Melhora as habilidades de resolução de problemas e a capacidade de abstração.



3

## Senso de organização

As habilidades que envolvem aprender a programar estão relacionadas com organização.

4

## Análise de Dados

Mesmo utilizando ferramentas que auxiliam no processo, códigos precisarão ser desenvolvidos.

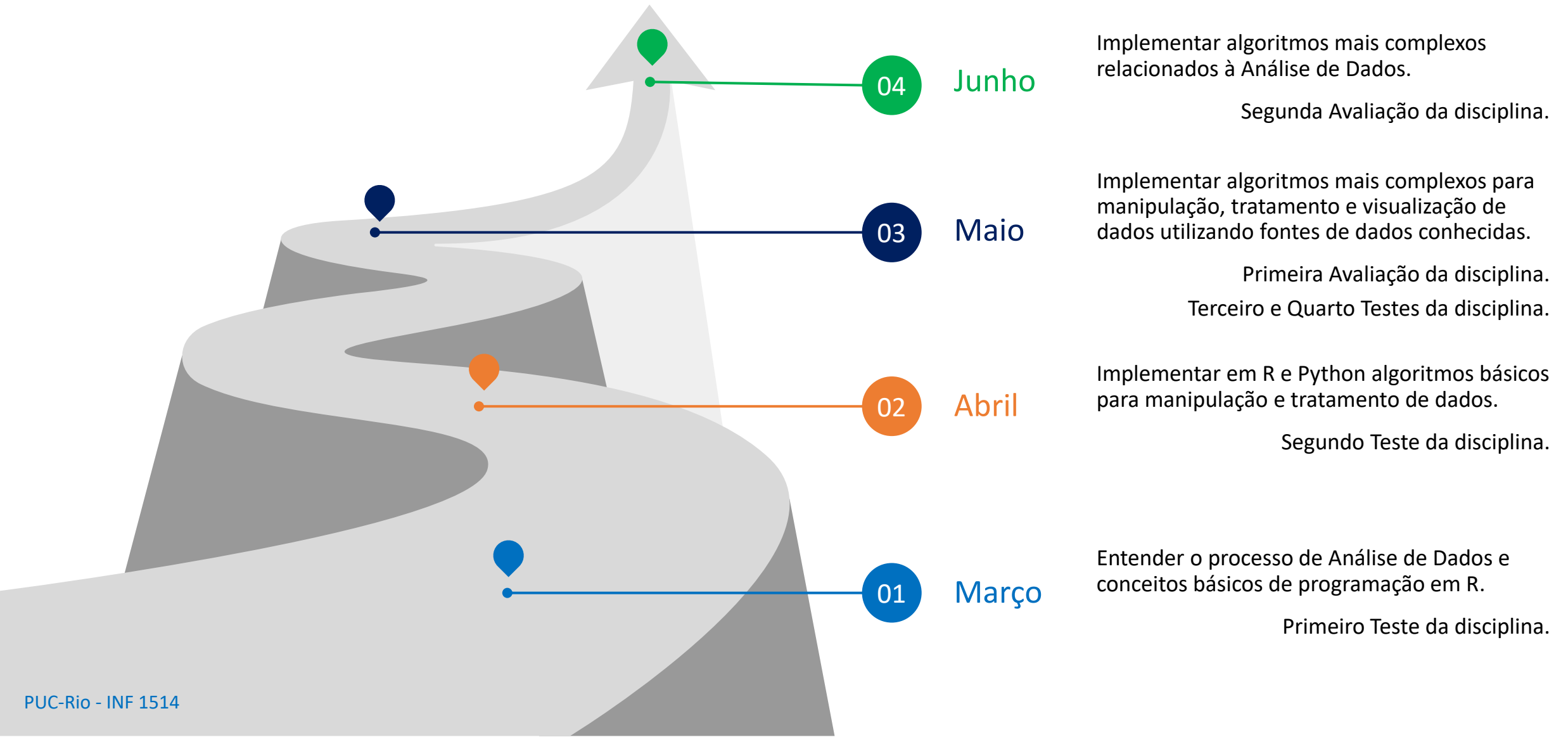
## Onde estamos

Obtendo uma visão geral do que é Análise de Dados e de sua importância, bem como a sua relação com o “saber programar”.

## Onde chegaremos

Ter a capacidade de implementar pequenos projetos de Análise de Dados utilizando as linguagens de programação R e Python.

# Nossa trilha



# O que é linguagem de programação de alto nível?

## Linguagens de programação de alto nível

- São compostas de símbolos mais complexos, inteligível pelo ser humano e não-executável diretamente pela máquina.

```
print("Bom dia!")
```

## Linguagens de programação de baixo nível

- Seus símbolos são uma representação direta do código de máquina.

```
section .data
    hello:      db 'Bom dia!',10
    helloLen:   equ $-hello
```

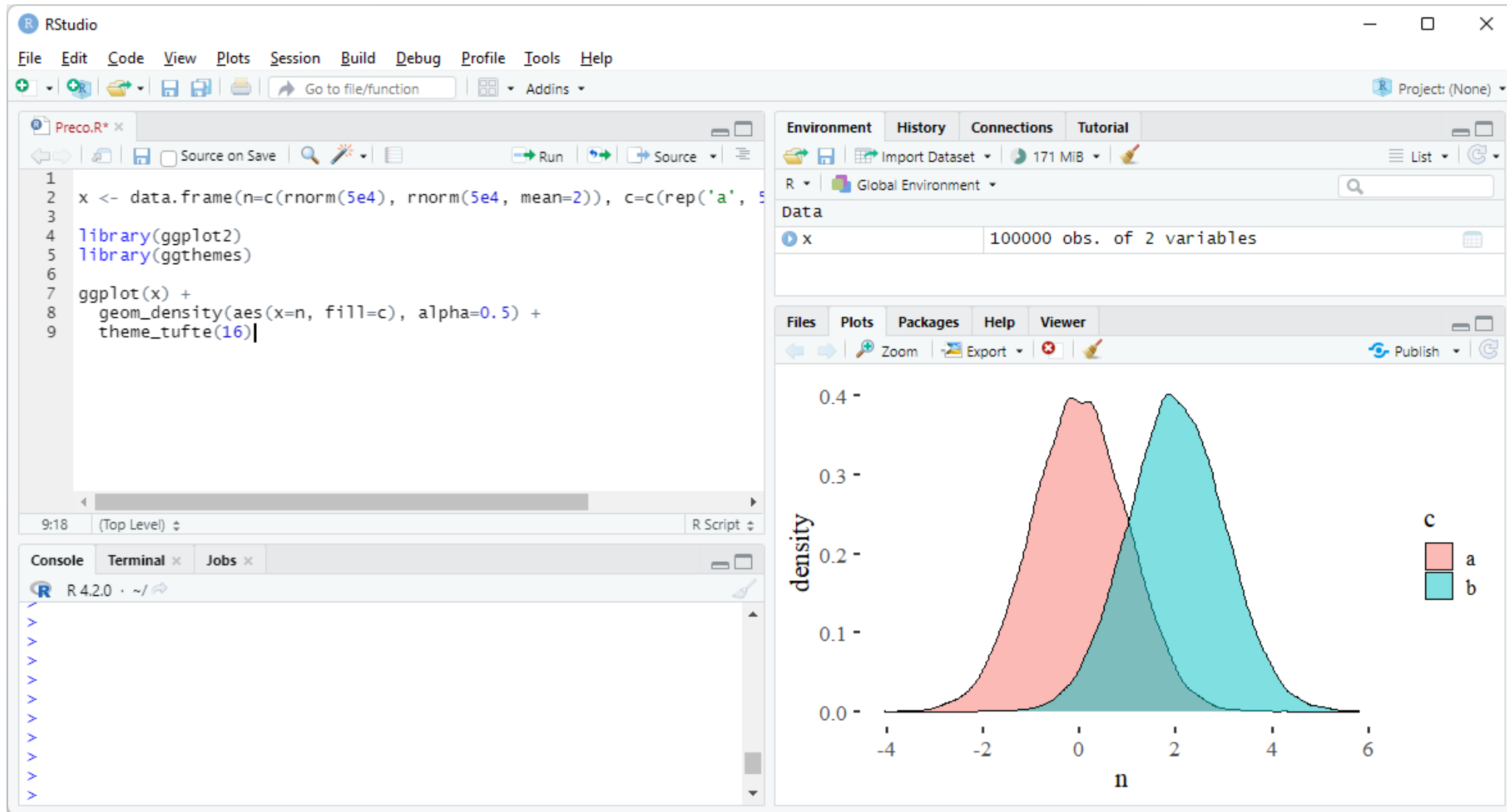
```
section .text
    global _start
```

```
_start:
    mov eax,4
    mov ebx,1
    mov ecx,hello
    mov edx,helloLen
    int 80h
    mov eax,1
    mov ebx,0
    int 80h;
```

# R

- R é uma **linguagem** e um **ambiente de desenvolvimento** integrado.
- Foi criada originalmente por Ross Ihaka e por Robert Gentleman na universidade de Auckland, Nova Zelândia, e foi desenvolvida por um **esforço colaborativo** de pessoas em vários locais do mundo.
- A sua estrutura de **código aberto** e de **software público e gratuito** atraiu um grande número de desenvolvedores, sendo que há hoje inúmeros pacotes para o R.
- O R disponibiliza uma **ampla variedade de técnicas estatísticas e gráficas**, incluindo modelagem linear e não linear, testes estatísticos clássicos, análise de séries temporais, classificação, agrupamento e outras.
- Tem uma **história longa e confiável** e uma **forte comunidade** de suporte no setor de dados.

# R Studio



# Python

- Foi lançada por Guido van Rossum em 1991 e, atualmente, possui um **modelo de desenvolvimento comunitário**, aberto e gerenciado pela organização sem fins lucrativos Python Software Foundation.
- É uma **linguagem de programação multiplataforma** que permite desenvolver aplicações para games, desktops, web e dispositivos móveis sendo também muito utilizada em aplicações que lidam com processamento de texto e machine learning.
- Combina uma sintaxe concisa e clara com os recursos poderosos de sua **biblioteca padrão** e por **módulos e frameworks** desenvolvidos por terceiros.
- A **facilidade de integração** com outras linguagens também é um fator de destaque.

# Spyder

The screenshot displays the Spyder IDE interface with the following components:

- Left Panel (Project Explorer):** Shows a tree view of the project structure. The 'Plots' folder is expanded, showing files like `plugin.py`, `chart_plot_example.py`, and `plugin.py - ipythonconsole`.
- Central Panel (Code Editor):** Displays the `plugin.py` file. The code defines a `Plots` class that inherits from `SpyderDockablePlugin`. It includes imports for `QtCore`, `Plugins`, `SpyderDockablePlugin`, `VariableExplorer`, and `Help`. The class has methods for `get_name`, `get_description`, `get_icon`, and `register`.
- Right Panel (Variable Explorer):** Shows a table of variables in the current scope. The table has columns for Name, Type, Size, and Value.
- Bottom Panel (Plots):** Displays two plots: a 3D surface plot and a polar plot.

Name	Type	Size	Value
bool	bool	1	True
data	Array of str128	(3, 3)	ndarray object of numpy module
datetime_object	datetime	1	2021-04-14 17:35:14.687085
df	DataFrame	(2, 2)	Column names: Col1, Col2
filename	str	53	/Users/Documents/spyder/spyder/tests/test_dont_use.py
li	list	5	['abcd', 745, 2.23, 'efgh', 70.2]
myset	set	3	{'2', '1', '3'}
r	float	1	6.46567886443
t	tuple	5	('abcd', 745, 2.23, 'efgh', 70.2)
tinylst	list	2	[123, 'efgh']
x	float64	1	1.1235123099439

conda: spyder-dev (Python 3.8.5) LSP Python: ready master Line 10, Col 1 UTF-8 LF RW Mem 64%



# Por que vamos ver R e Python?

- De acordo com as características do projeto uma “linguagem” pode ser **mais indicada do que a outra** ou mesmo podem ser **utilizadas em conjunto** aproveitando-se o melhor de cada uma.
- **Traçar paralelos** (prós e contras) entre linguagens demonstrando de forma prática diferentes abordagens, conceitos e recursos que podem ser encontrados em outros ambientes e linguagens.
- Apesar de existirem muitas outras possibilidades, estas duas linguagens tem **polarizado as discussões** sobre que ferramenta utilizar para Análise de Dados.
- As duas linguagens **são simples** (e gratuitas) para instalar e relativamente fáceis de começar a usar servindo como excelente ponto de partida para se aprender a programar para Análise de Dados.
- Vamos **começar por R** por apresentar menor curva de aprendizado.

## Exemplo em R



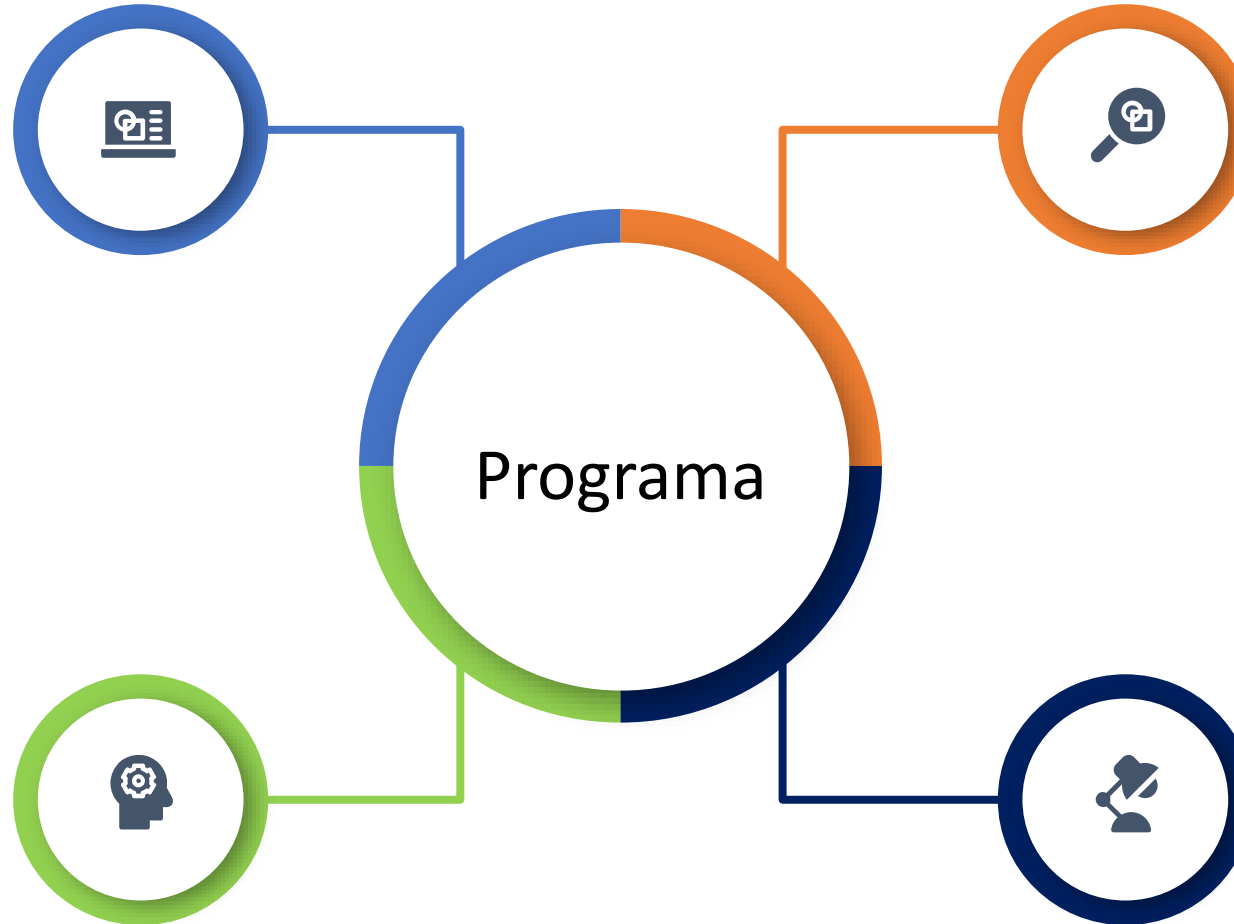
# INF 1514

## Base de programação

Conceitos básicos;  
importação e exportação  
de dados; repetição;  
condicionais; funções.

## Manipulação de dados

Subconjuntos;  
ordenação;  
transformações; merge;  
fontes de dados.

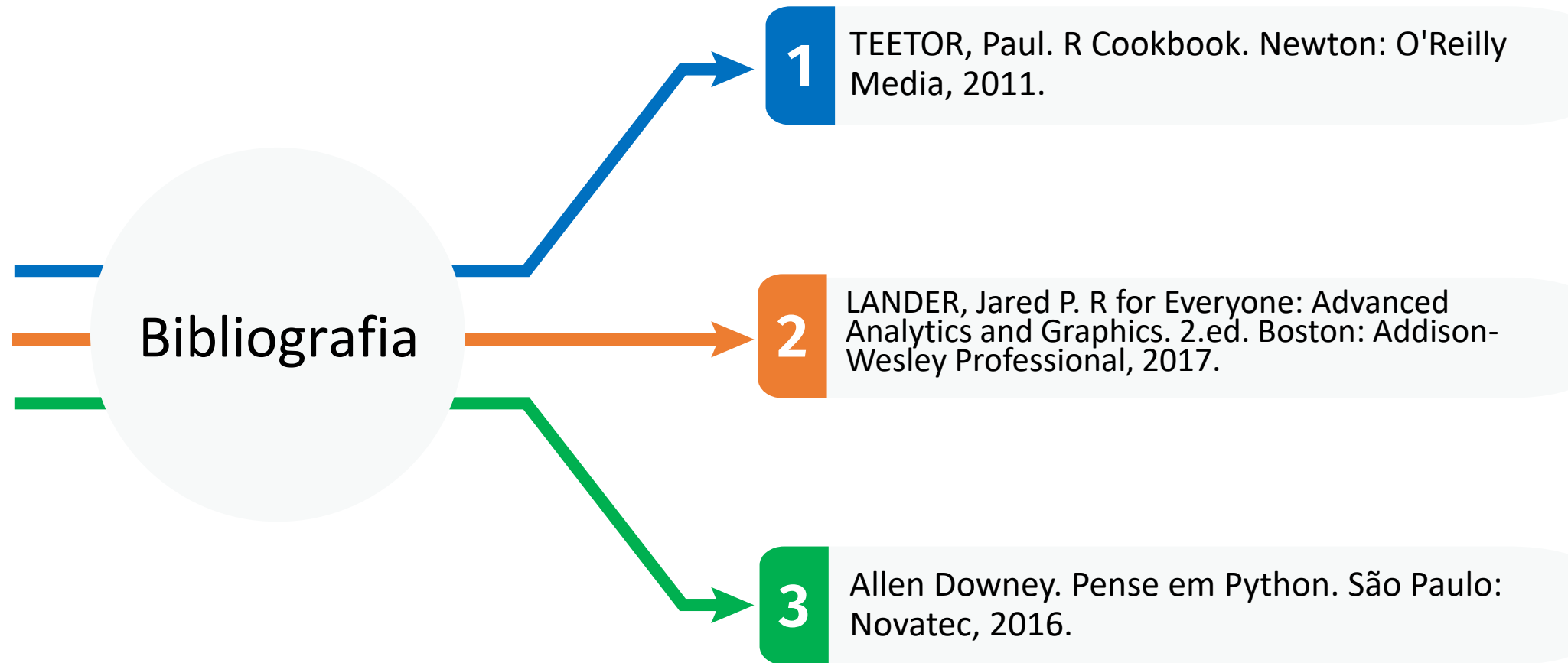


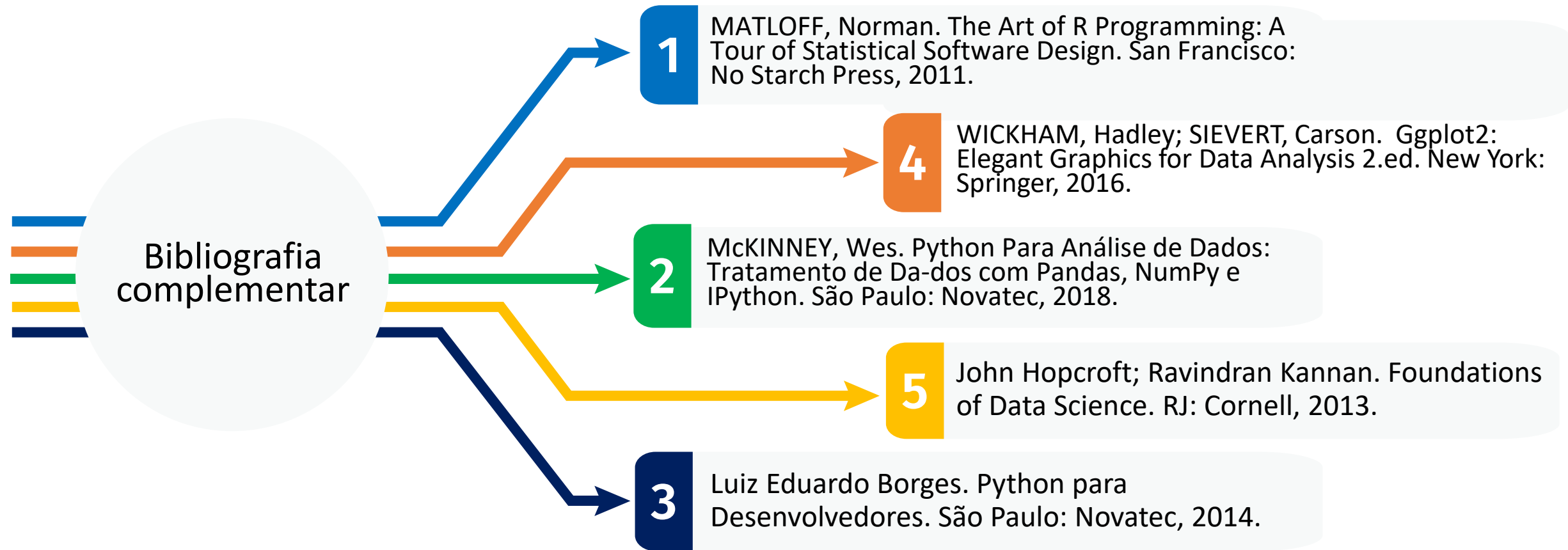
## Visualização de dados

Dados temporais;  
espaciais; espaço-  
temporais.

## Outros tópicos

Web scraping;  
organização de código;  
documentação;  
versionamento de código.





# Critério de Avaliação

- Critério 3.
- Nota da G1:
  - Teste 1, Teste 2 e Prova.
  - $G1 = \left( \frac{Teste\ 1 + Teste\ 2}{2} \right) * 0,4 + Prova * 0,6$
- Nota da G2:
  - Teste 1, Teste 2 e Trabalho.
  - $G1 = \left( \frac{Teste\ 1 + Teste\ 2}{2} \right) * 0,4 + Trabalho * 0,6$
- Média Final (MF):
  - $MF = \left( \frac{G1 + G2}{2} \right)$
  - $MF \geq 5$ , aprovado.