

Data Science

com *software open source*+livre+grátis:
um bom casamento

Ricardo São João

IPSantarém, CEAUL, CIDNUR, CEG-UAb

ricardo.sjoao@esg.ipsantarem.pt

Outubro, 2025

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

Documentos
dinâmicos

Requisitos

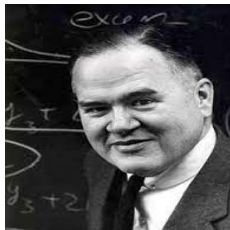
Aplicações

- 1 *Data Science*
- 2 *Software livre*
- 3 *Software open source*
- 4 Gráficos
- 5 Documentos dinâmicos
- 6 Requisitos
- 7 Aplicações

O que é ?

Não existe uma definição formal e consensual relativamente a *Data Science*/Ciência dos Dados.

Já em 1962 (há 63 anos) o matemático/estatístico John Tukey aborda a questão na publicação intitulada *The future of data analysis*¹



fonte:<https://4.bp.blogspot.com>

¹Tukey, J. W. (1962). The future of data analysis. The annals of mathematical statistics, 33(1), 1-67.

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

❖ *Data Science* é referido como um campo multidisciplinar (matemática, estatística, informática, engenharia) onde são utilizadas metodologias distintas com o objetivo de **extrair valor (informação) dos dados**.

❖ Nas últimas duas décadas a Ciência dos Dados tem assistido a um forte crescimento e protagonismo sendo os dados atualmente apelidados de “**o novo petróleo**”.

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

★ A profissão “Cientista de Dados” é das mais procuradas e melhor remuneradas.

★ Ocupou o 6º lugar das 10 profissões mais procuradas em 2024.
(Fonte: [Linkedin](#))

[https://www.linkedin.com/pulse/
future-proof-your-career-top-10-in-demand-jobs-xlsvc](https://www.linkedin.com/pulse/future-proof-your-career-top-10-in-demand-jobs-xlsvc)

Alguns números que dão que pensar...

- ✿ o volume de dados na internet no final de 2020 foi estimado em 44 zettabytes (1 zettabyte= 1.000.000.000.000.000.000= 10^{21});
- ✿ até 2025, espera-se que a quantidade de dados gerados a cada dia atinja 463 exabytes (1 exabyte= 10^{18});

²fonte: <https://seedscientific.com>

Alguns números que dão que pensar...

✿ o volume de dados na internet no final de 2020 foi estimado em 44 zettabytes (1 zettabyte= 1.000.000.000.000.000.000= 10^{21});

✿ até 2025, espera-se que a quantidade de dados gerados a cada dia atinja 463 exabytes (1 exabyte= 10^{18});

entretanto já me “perdi” com tantos zeros.

²fonte: <https://seedscientific.com>

Alguns números que dão que pensar...

✿ o volume de dados na internet no final de 2020 foi estimado em 44 zettabytes (1 zettabyte = 1.000.000.000.000.000.000 = 10^{21});

✿ até 2025, espera-se que a quantidade de dados gerados a cada dia atinja 463 exabytes (1 exabyte = 10^{18});

entretanto já me “perdi” com tantos zeros.

✿ *Google, Facebook, Microsoft e Amazon* armazenam pelo menos 1.200 petabytes de informação (1 petabyte = 10^{15});

✿ a cada minuto são gastos mais de 12 mil milhões de dólares em compras na WWW;

✿ em 2030, nove em cada dez pessoas com idade ≥ 6 anos será digitalmente ativa.

²fonte: <https://seedscientific.com>

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Porém a análise de tal volume de informação exige recursos/meios, onde o *software* seguramente é um deles.

A escolha pode recair em *software* comercial ou ...

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Porém a análise de tal volume de informação exige recursos/meios, onde o *software* seguramente é um deles.

A escolha pode recair em *software* comercial ou ...

- livre;

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Porém a análise de tal volume de informação exige recursos/meios, onde o *software* seguramente é um deles.

A escolha pode recair em *software* comercial ou ...

- livre;
- *open source*;

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Porém a análise de tal volume de informação exige recursos/meios, onde o *software* seguramente é um deles.

A escolha pode recair em *software* comercial ou ...

- livre;
- *open source*;
- gratuito.

Software livre: o quê significa ?

O termo “livre” não está associado à ideia de não pagamento na aquisição de *software*.



fonte:<https://foreignpolicyi.org>



fonte:<https://poupaeganha.pt>

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Não se trata de gratuidade (“de graça” ou “de borla”).

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Não se trata de gratuidade (“de graça” ou “de borla”).

Estaremos então a falar de
que conceito ?

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações



fonte:<http://apoesiadaisamar.blogspot.com/>

Significa que os seus
utilizadores possuem
“liberdade” (são livres).

Ricardo São
João

Data Science

Software livre

Software open
source

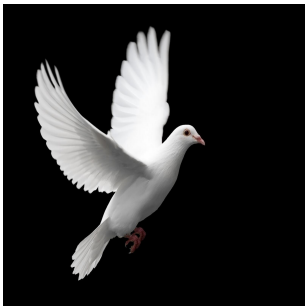
Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações



fonte: <http://apoesiadaisamar.blogspot.com/>

Significa que os seus
utilizadores possuem
“liberdade” (são livres).

A noção de liberdade assenta em
quatro pilares essenciais:

- liberdade na execução do *software*/programa;
- liberdade no desenvolvimento e alteração do código fonte;
- liberdade na redistribuição de cópias;
- liberdade na redistribuição de versões modificadas.

Ricardo São
João

Data Science

Software livre

**Software open
source**

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Origem

Nasce como alternativa ao código proprietário presente em *software* comercial (pagamento de licenças).

No *software open source*, o(s) seu(s) autor(es) abdica(m) da propriedade intelectual do código de forma a que outros utilizadores possam tirar benefício.

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Características:

- Usualmente comunga dos quatro pilares do *software* livre;
- Não aplica qualquer tipo de discriminação/restricção aos seus utilizadores.

Vantagens:

- *Download* acessível e gratuito;
- Trabalho em plataformas colaborativas podendo usufruir dos contributos de outros utilizadores;
- Beneficia de uma melhoria constante impulsionada pela comunidade de utilizadores*;
- Redução de custos.

* Tal realidade não é expectável num *software* de código proprietário.

Ricardo São
João

Data Science

Software livre

Software open
source

**Características &
Vantagens**

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Analogia *open source* + livre com as maçãs:

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

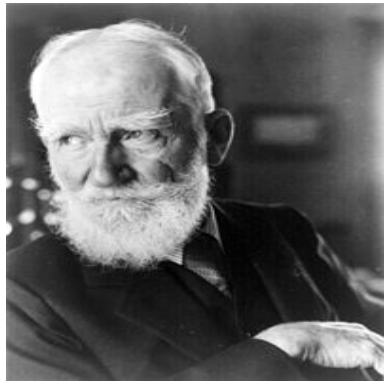
Documentos
dinâmicos

Requisitos

Aplicações

Analogia open source + livre com as maçãs:

“Se tu tiveres uma maçã e eu tiver uma maçã, trocando essas maçãs, continuaremos, **cada um, a ter uma maçã.** Mas se tu tiveres uma ideia e eu tiver uma ideia, trocando essas ideias, **cada um de nós passará a ter duas ideias.**”



George Bernard Shaw (1856-1950)
Nóbel da literatura em 1925;
Óscar de melhor argumento
adaptado - 1939.

R: *open source* + livre + grátis um bom “casamento”!

- O R é um projecto *open source*, que utiliza a linguagem de programação R, linguagem por excelência no tratamento e análise de dados;

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

R: *open source* + livre + grátis um bom “casamento”!

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

- O R é um projecto *open source*, que utiliza a linguagem de programação R, linguagem por excelência no tratamento e análise de dados;
- Surge em 1993 e foi criado originalmente por **Ross Ihaka** e por **Robert Gentleman** no departamento de Estatística da Universidade de Auckland, Nova Zelândia.

R: *open source* + livre + grátis um bom “casamento”!

Ricardo São
João

Data Science

Software livre

Software *open
source*

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

- O R é um projecto *open source*, que utiliza a linguagem de programação R, linguagem por excelência no tratamento e análise de dados;
- Surge em 1993 e foi criado originalmente por **Ross Ihaka** e por **Robert Gentleman** no departamento de Estatística da Universidade de Auckland, Nova Zelândia.

- Funciona em todos os sistemas operativos e tem mais de 22 mil pacotes gratuitos transversais a todas as ciências.



Fonte:

<https://cran.r-project.org/>

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

A promotional banner for GitHub. On the right, a large blue globe is composed of a grid of dots, with several red lines representing code paths or connections. In the bottom right corner, a small cartoon astronaut in a blue and white suit is looking up at the globe. The background is a dark blue gradient.

**Where the world
builds software**

Millions of developers and companies build, ship, and maintain their software on GitHub—the largest and most advanced development platform in the world.

Email [Sign up for GitHub](#)

56+ million Developers	3+ million Organizations	100+ million Repositories	72% Fortune 50
---------------------------	-----------------------------	------------------------------	-------------------

Fonte: <https://github.com/>

Ricardo São
João

Data Science

Software livre

Software open
source

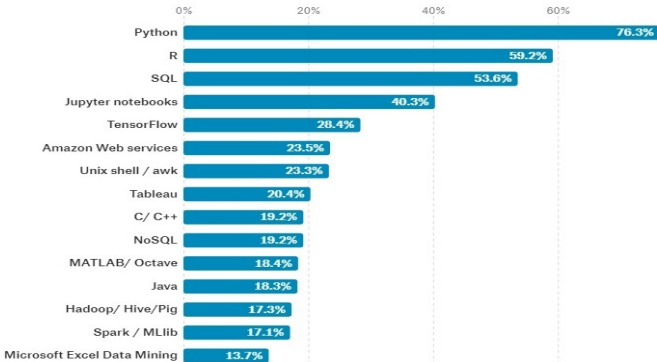
Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações



Fonte: <https://stackoverflow.blog>

R: Educação, Investigação, Data Science, Business Analyst, Data Analyst, Data Miner, Operations Researcher, Predictive Modeler, Marketing Researcher, Jornalismo, SIG's, . . .

Ricardo São João

Data Science

Software livre

Software open source

Características & Vantagens

Gráficos

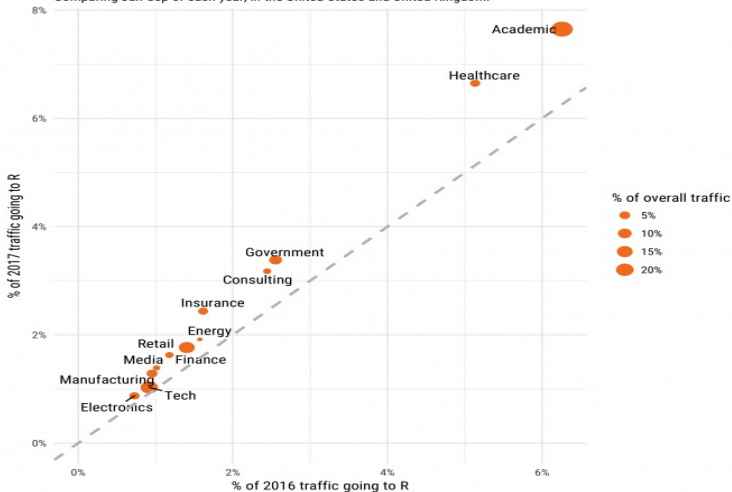
Documentos dinâmicos

Requisitos

Aplicações

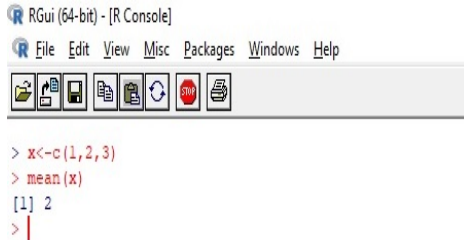
Traffic by industry to R

Comparing Jan-Sep of each year, in the United States and United Kingdom.



Fonte: <https://blog.revolutionanalytics.com/popularity/>

As instruções no R são dadas por comandos em linhas de código numa consola o que “inibe” à partida os seus (futuros) utilizadores.



no entanto . . .

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

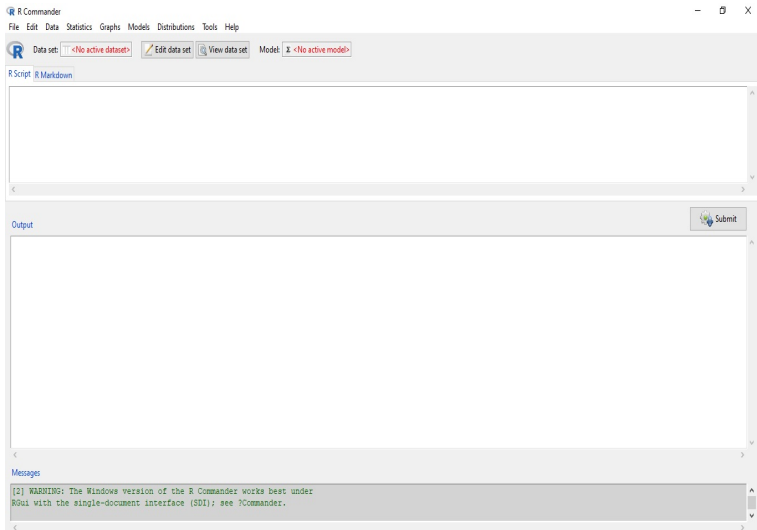
Aplicações

...existem interfaces gráficas (*Graphical User Interface - GUI*) que permitem uma utilização mais fácil e intuitiva do R.

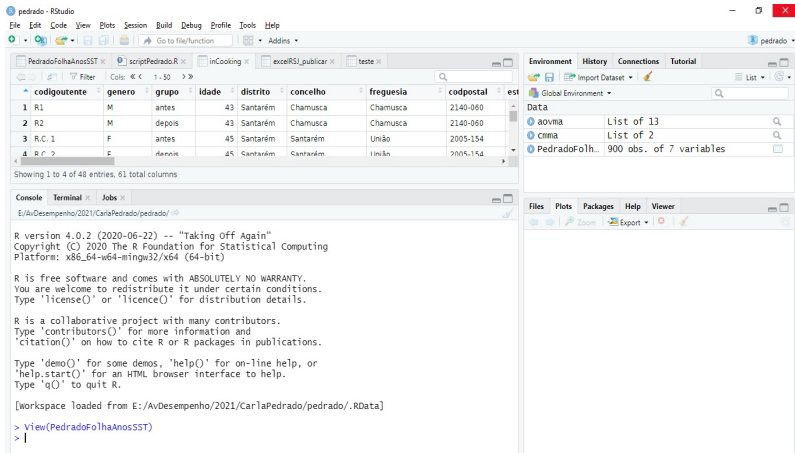
Dentre as várias GUI's disponíveis destacam-se duas:

- R Commander;
- RStudio.

Ambas as GUI necessitam ter o R instalado acessível em
<https://cran.r-project.org/>



<https://cran.r-project.org/web/packages/Rcmdr/index.html>



The screenshot shows the RStudio interface with the following components:

- Environment Pane:** Shows a 'Global Environment' with a 'Data' section containing:
 - aovma: List of 13
 - cnma: List of 2
 - PedradoFolh.: 900 obs. of 7 variables
- Table View:** Displays a data table with columns: genero, grupo, idade, distrito, concelho, freguesia, and codpostal. The first four rows are:

	genero	grupo	idade	distrito	concelho	freguesia	codpostal
1	R1	M	antes	43	Santarém	Chamusca	2140-060
2	R2	M	depois	43	Santarém	Chamusca	2140-060
3	R.C. 1	F	antes	45	Santarém	Santarém	2005-154
4	R.C. 2	F	depois	45	Santarém	Santarém	2005-154
- Console:** Shows the R version (4.0.2), copyright information, and standard startup messages. The last command executed is `> View(PedradoFolhaAnosSST)`.

<https://rstudio.com/>

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

O R dispõem de *plug-ins* (módulos) que permitem elaborar análises e gráficos específicos para diferentes áreas de estudo.

Pacotes (*packages*) e Bibliotecas (*libraries*)

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações



Sempre que necessário o **R** permite o *download* de forma gratuita de mais de 22 mil pacotes !

```
install.packages("nomepacote");library(nomepacote)
```

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

<https://r-graph-gallery.com/index.html>

✿ Vamos ver/implementar agora alguns ...

<https://r-graph-gallery.com/9-ordered-boxplot.html#grouped>

[//r-graph-gallery.com/9-ordered-boxplot.html#grouped](https://r-graph-gallery.com/9-ordered-boxplot.html#grouped)

[https://r-graph-gallery.com/](https://r-graph-gallery.com/199-correlation-matrix-with-ggally.html)

[199-correlation-matrix-with-ggally.html](https://r-graph-gallery.com/199-correlation-matrix-with-ggally.html)

<https://r-graph-gallery.com/44-polynomial-curve-fitting.html>

[//r-graph-gallery.com/44-polynomial-curve-fitting.html](https://r-graph-gallery.com/44-polynomial-curve-fitting.html)

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

**Documentos
dinâmicos**

Requisitos

Aplicações

Quantas vezes ...

- não teve de reformular um documento/relatório com base em novas informações/dados ?

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

**Documentos
dinâmicos**

Requisitos

Aplicações

Quantas vezes ...

- não teve de reformular um documento/relatório com base em novas informações/dados ?
- a deteção de um valor/parâmetro incorreto numa análise, não comprometeu o estudo realizado?

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

**Documentos
dinâmicos**

Requisitos

Aplicações

Quantas vezes ...

- não teve de reformular um documento/relatório com base em novas informações/dados ?
- a deteção de um valor/parâmetro incorreto numa análise, não comprometeu o estudo realizado?
- não gostaria de poder reproduzir os mesmos valores/resultados com base na mesma informação de um artigo científico ou relatório técnico ?

✿ Vamos criar um relatório com o rmarkdown.



<https://rmarkdown.rstudio.com/>

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

Documentos
dinâmicos

Requisitos

Aplicações

Ricardo São
João

Data Science

Software livre

*Software open
source*

*Características &
Vantagens*

Gráficos

**Documentos
dinâmicos**

Requisitos

Aplicações



<https://shiny.rstudio.com/>

Ricardo São
João

Data Science

Software livre

Software open
source

Características &
Vantagens

Gráficos

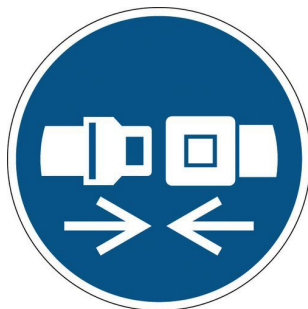
Documentos
dinâmicos

Requisitos

Aplicações

- 1 Instale o **R** a partir do CRAN (*The Comprehensive R Archive Network*) acessível em <https://cran.r-project.org/>
- 2 Instale o **RStudio** acessível em <https://posit.co/download/rstudio-desktop/>

✿ Vamos começar o nosso *tour* com alguns exemplos.
Preparados . . . apertem os vossos cintos de segurança !



fonte imagem:<https://www.kaiserkraft.pt>

Ricardo São João

Data Science

Software livre

Software open source

Características & Vantagens

Gráficos

Documentos dinâmicos

Requisitos

Aplicações

Hitler



Mussolini



Churchill



Eisenhower



Stalin



Franco



Charles de Gaulle



Tito



imagens retiradas de <https://pt.wikipedia.org>

O modelo MDS permite a representação espacial de **perceções** e **preferências** de consumidores/inquiridos ao nível da (dis)similaridade entre objetos

Foram avaliadas as parecenças/semelhanças ³ entre alguns políticos na II Grande Guerra Mundial com base numa escala ordinal, cuja informação apresenta-se na seguinte matriz:

> `politicosIIguerramundial`

	Hitler	Mussolini	Churchill	Eisenhower	Stalin	Attlee	Franco	De_Gaulle	Mao_Tse	Truman	Chamberlain	Tiro
Hitler	0											
Mussolini	5	0										
Churchill	11	14	0									
Eisenhower	15	16	7	0								
Stalin	8	13	11	16	0							
Attlee	17	18	11	16	15	0						
Franco	5	3	12	14	13	16	0					
De_Gaulle	10	11	5	8	11	12	9	0				
Mao_Tse	16	18	16	17	12	16	17	13	0			
Truman	17	18	8	6	14	12	16	9	12	0		
Chamberlain	12	14	10	7	16	9	10	11	17	9	0	
Tiro	16	17	8	12	12	13	12	7	10	11	15	0

O modelo MDS não métrico irá ajudar a interpretar a informação.

³adaptado de https://rpubs.com/Hasantha_APS_1701/

Ricardo São João

Data Science

Software livre

Software open source

Características & Vantagens

Gráficos

Documentos dinâmicos

Requisitos

Aplicações



No eixo das abcissas podemos ter em conta a **dimensão ideologia**. À medida que caminhamos para a esquerda do eixo das abcissas a ideologia nazi ganha relevo. No eixo y podemos ter a **dimensão localização geográfica**. Políticos posicionados mais acima do eixo correspondem a países orientais.

A título ilustrativo vamos considerar um índice constituído por 10 itens avaliados numa escala tipo likert com as seguintes pontuações: 0 (no); 5 (sometimes) e 10 (yes).

Item	Description
1	Do you have difficult opening your mouth wide?
2	Do you have difficulty moving your jaw to the sides?
3	Do you feel fatigue or muscle pain when you chew?
4	Do you have frequent headaches?
5	Do you have neck pain or a stiff neck?
6	Do you have ear aches or pain in that area?
7	Have you ever noticed any noise while chewing or opening your mouth?
8	Do you have any habits such as clenching or grinding your teeth?
9	Do you feel that your teeth do not come together well?
10	Do you consider yourself a nervous person?

No **R** é possível realizar uma análise detalhada. Ora vejamos ...

Análise bivariada em tabelas de contingência

Considere um estudo constituído por 100 leitores de poetas portugueses. A seguinte tabela de contingência mostra a frequência dos leitores tendo em conta o seu sexo e o poeta presentemente lido.

Poeta	Sexo do leitor	
	M	F
Luís de Camões (1524-1580)	12	26
Fernando Pessoa (1888-1935)	8	10
Sophia Andresen (1919-2004)	30	14



Existirá alguma relação entre o sexo do leitor e a escolha do poeta ?
 Se sim, com que intensidade ? Veja uma abordagem no **R** ...

Exemplo

A base de dados penguins disponível no pacote palmerpenguins do R retrata as medidas de um conjunto de 344 pinguins de três espécies (Adelie, Chinstrap and Gentoo) que passam pelo Arquipélago Palmer. Através da ANOVA procure dar resposta à seguinte questão: “O comprimento médio das barbatanas apresenta diferenças estatisticamente significativas nas 3 espécies de pinguins?”^a

^a adaptado de <https://statsandr.com>. Imagens retiradas de: penguinlovers.de; 4.bp.blogspot.com; images.fineartamerica.com



Adelie



Chinstrap



Gentoo

O presente exemplo consiste na informação recolhida em 1974 pela revista *Motor Trend US magazine*^a e diz respeito ao ensaio de 32 modelos de automóveis, comercializados entre os anos de 1973 e 1974, onde são avaliadas características técnicas e o consumo. A base de dados está presente no **R** acessível com a instrução `data(mtcars)`. Para maiores detalhes faça `?mtcars`

^aainda comercializada <https://www.motortrend.com/>



Objetivo

Com base num modelo de RL procure explicar o consumo de um automóvel com base em algumas especificações técnicas.

Ricardo São João

Data Science

Software livre

Software open source

Características & Vantagens

Gráficos

Documentos dinâmicos

Requisitos

Aplicações

Obrigado pela Vossa **PPA** !
Presença, **P**aciência e **A**tenção.

Ricardo São João

ricardo.sjoao@esg.ipsantarem.pt