

Introdução ao R: aplicações na linguística com corpos

Diana Santos

EBraLC 2012: Material para ler e aplicar antes do curso

1 Primeiros passos: Entrar e sair

Carregue na imagem do R, ou invoque da linha de comandos o comando “R”.

O interpretador de R tem em geral o “pronto” `>`. Aí pode começar a escrever comandos.

O comando mais importante é como sair do interpretador:

```
> q()
```

O segundo comando mais importante é pedir ajuda

```
> help()
```

Também pode ser interessante e importante saber onde está, ou seja, como é que o R interpreta o ambiente em que está instalado:

```
> getwd()
```

Conforme o sistema operacional em que trabalha, poderá obter respostas como

```
[1] "e:/work/dssantos"  
[1] "/hf/sokrates/ilos-u1/dssantos"
```

Presumo que saberá interpretá-las, e se quiser que a diretoria de trabalho com o R seja outra, poderá instruir o R para modificar essa variável:

```
> setwd("m:Rwork/")
```

ou

```
> setwd("/linguateca/cursos/R/")
```

Para se familiarizar com a forma como o R funciona e responde, pode experimentar alguns comandos do tipo “máquina de calcular”

```
> 3+4
[1] 7
> "olá"
[1] "olá"
> c("23","1","batatas","cenouras")
[1] "23"      "1"      "batatas" "cenouras"
> bobagem <- c("23","1","batatas","cenouras")
> bobagem
[1] "23"      "1"      "batatas" "cenouras"
> algo
Error: object 'algo' not found
> 'desculpe'
[1] "desculpe"
> mas()
Error: could not find function "mas"
> help
function (topic, package = NULL, lib.loc = NULL, verbose = getOption("verbose"),
  try.all.packages = getOption("help.try.all.packages"), help_type = getOption("help.type"))
{
  types <- c("text", "html", "pdf")
  if (!missing(package))
    ...
}
>
```

Mais algumas explorações aritméticas

```
> sqrt(49)
[1] 7
> 24:100
[1] 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42
[20] 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61
[39] 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80
[58] 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99
[77] 100
```

Isso dá para compreender, espero, um conjunto de realidades básicas sobre o R, nomeadamente, que

- cada linha de resposta vem numerada, começando com 1, e a ordem do primeiro caso dessa linha é sempre impressa no lado esquerdo, entre []

- quando são dadas instruções, o R executa o comando correspondente e devolve o resultado
- quando se manda simplesmente um nome, o R vai ver se a esse nome já foi dado algum sentido. Se foi, dá-nos o valor atribuído (que até pode ser uma descrição de código), senão dá uma mensagem de erro
- quando se lhe diz um valor, ele ecoa esse valor
- há diferença entre nomes de variáveis, e cadeias de caracteres (marcadas por aspas ou plicas)
- as funções no R são designadas por NOME(argumentos), quando não têm argumentos, como a função `help`, têm de ser marcadas por `()` a seguir ao nome

2 Comunicar com o exterior

Se não quisermos simplesmente brincar dentro do R, ou usá-lo como um oráculo ou máquina de calcular, mas sim interagir com dados pré-existentes e produzir novos dados ou figuras, é preciso saber como ler e escrever arquivos e figuras. A maneira mais fácil para começar é ler de um ficheiro existente na rede, por isso os primeiros exercícios irão proceder de um conjunto de dados que pus na minha página. Mas depois vão progressivamente usando arquivos dentro do vosso computador.

```
> CondivSimples <- read.table("http://folk.uio.no/dssantos/cursoR/condiv.txt")
> CondivSimples
      V1      V2      V3      V4      V5
1  dataset colour  theme decade variety
2  fashion50BR  1512  fashion     50     BR
3  fashion50PT  2351  fashion     50     PT
[...]
```

Ao inspecionar o conteúdo do que lemos é imediatamente aparente que a primeira linha é especial, porque é um cabeçalho. Existe uma forma simples de dizer isso ao R:

```
> CondivSimples <- read.table("http://folk.uio.no/dssantos/cursoR/condiv.txt", h
> CondivSimples
      dataset colour  theme decade variety
1  fashion50BR  1512  fashion     50     BR
2  fashion50PT  2351  fashion     50     PT
[...]
```

O conteúdo deste arquivo é numa forma que é a mais usada para fazer análises estatísticas, tendo uma observação por linha e o valor de várias variáveis, cada variável correspondendo a uma coluna. Em inglês, “dataframe”, proponho chamar em português uma “folha de registo”.

Uma folha de registo é como uma matriz em que tanto às linhas como às colunas se podem dar nomes, a maior diferença é que um matriz pressupõe que todos os elementos são do mesmo tipo, enquanto numa folha de registo não.

Uma forma de obter rapidamente as dimensões (número de linhas e de colunas) é usar a função `dim`:

```
> dim(CondivSimples)
[1] 18  5
```

Para obter apenas uma coluna, pode-se simplesmente usar o nome da folha de registo seguida de cifrão e do nome da coluna:

```
CondivSimples$colour
[1] 1512 2351 1342 6481 1525  669  719  694  468  318 1167  614  136  583  127
[16]  993   97  265
```

Também podemos ver que a primeira coluna é sempre diferente, e que presumivelmente foi usada para identificar de forma única cada linha. Por isso, é possível informar o R disso também, através da indicação `rownames=número da coluna`. Isto não é contudo tão comum, porque as folhas de registo costumam ter milhares de linhas e nem sempre faz sentido dar-lhes um identificador único.

```
> CondivSimples <- read.table("condiv.txt",header=TRUE,row.names=1)
> CondivSimples
      colour  theme decade variety
fashion50BR   1512 fashion     50     BR
fashion50PT   2351 fashion     50     PT
[...]
```

Isso implica que agora existem apenas quatro colunas na folha de registo, e não cinco como anteriormente.

```
> dim(CondivSimples)
[1] 18  4
```

Podemos estar interessados em criar uma nova folha de registo só com uma parte das informações de outra, para isso existe a função `subset`. No exemplo seguinte, selecionamos apenas as linhas relativas ao futebol e as colunas com cor e variante:

```
> SoFutebol<-subset(CondivSimples,theme=="football",c(1,4))
> SoFutebol
```

	colour	variety
football150BR	719	BR
football150PT	694	PT
football170BR	468	BR
football170PT	318	PT
football2000BR	1167	BR
football2000PT	614	PT

Vamos agora escrever esta folha de registo na vossa diretoria de trabalho.

```
> write.table(SoFutebol, file="futebol.txt")
```

Além de poder ir ver, através do seu sistema operativo, o que está nesse arquivo, outra forma é tornar a lê-lo para o R:

```
read.table("futebol.txt")
```

Vamos agora criar uma tabela à mão, e depois escrevê-la num arquivo externo. Embora esta não seja uma maneira natural de trabalhar, é uma maneira adequada de nos familiarizarmos com o R.

3 Figuras e gráficos

Naturalmente, até agora o R não se destacou de qualquer outra linguagem de programação e/ou de manipulação de texto, de matrizes ou de outras estruturas de dados.

A questão das capacidades gráficas é talvez o primeiro caso em que os usuários se apercebem da diferença e vantagem em usar o R.

Em primeiro lugar, é preciso compreender que o R dispõe de várias janelas, e que por omissão abre uma janela para todos os gráficos. Ou seja, ao desenharmos um gráfico ou figura, deixamos de ver e de poder usar a anterior. A não ser que a guardemos, ou que explicitamente indiquemos que queremos uma nova janela.

```
> plot(colour ~ variety, data=CondivSimples)
> hist(CondivSimples$colour)
> plot(CondivSimples$colour)
```

Se estivermos a trabalhar com uma dada folha de registo, dá jeito não ter sempre que estar a escrever o nome todo, e usar apenas os nomes associados a ela. Isso consegue-se com o comando `attach`:

```
> attach(CondivSimples)
```

Refaça as figuras anteriores depois de se ter associado à folha de registo, para ver como é mais simples.

Uma das coisas que convém compreender em relação ao R é que tem um mecanismo de simplificação (ou omissão) que permite várias formas diferentes de invocar uma “mesma” função, especificando ou não um número variável de coisas e de argumentos. (Tal com em linguagem natural um verbo n-ário pode ser usado com menos argumentos...) Assim, a função `plot` pode dar resultados, ou melhor gráficos, completamente diferentes conforme os argumentos e seus tipos. Além disso podem especificar-se muitíssimas escolhas diferentes, como por exemplo o título da figura, a forma e a cor dos pontos, o nome das abcissas e das ordenadas, a escala, etc etc etc. Alguns exemplos simples – mas antes disso, quero chamar a atenção para o comportamento do R quando os comandos são maiores do que uma linha: o R fica à espera do resto, e mostra que está à espera, escrevendo `+` em vez de `>`:

```
> plot(colour, col="green", pch=3, type="b")
> plot(colour, col="red", pch=2, las="1", tcl="0.2", xlab="casos",
+ ylab="Quantas palavras de cor", type="h")
> plot(colour, col="red", pch=2, las="1", xlab="casos",
+ legend(4,6481,"Máximo"))
```

Mas para explorar bem esta vertente, nada como fazer `help(plot)` e ler as várias opções acessíveis, ou consultar o cartão de referência do R[15].

De qualquer maneira, é importante explicar como obter as figuras num formato que dê para incluir noutras aplicações, quer artigos, quer apresentações. Isso consegue-se simplesmente direcionando o lugar para onde se imprime através de comandos com o nome do formato que se pretende. Se quiser um pdf, basta fazer

```
> pdf("simple.pdf")
> plot(colour~theme)
> dev.list()
```

Se quiser voltar à sua janela, basta ver quais as janelas através do comando `dev.list`, e depois selecionar a que quer através do respetivo número

```
> dev.list()
X11cairo      pdf
      2        3
> dev.set(2)
```

Enquanto não fechar a janela em que está a escrever, não pode geralmente abri-la com outros programas, por isso é importante que, quando estiver pronta, feche com o comando `dev.off()`, que pode usar argumentos, ou que feche a última janela aberta.

```
> dev.off()
X11cairo
      2
```

O mesmo se passa para formatos como `png`, `bmp`, etc.

Se quiser abrir várias janelas no seu computador ao mesmo tempo, e não escrever por cima, basta usar `dev.new`.

4 Preservar o trabalho entre sessões

Embora tal não seja apropriado para principiantes, dá muito jeito pder preservar o trabalho entre sessões quando se trabalha intensivamente e temos de interromper no meio de um percurso cheio de ingredientes e tarefas feitas.

Nesses casos, pode-se instruir o R para guardar tudo para quando nós voltarmos. O R tem um formato especial, do tipo `.Rdata` e comandos para guardar e recuperar o que se guardou.

```
> save.image("tudo.Rdata")
```

Saia do R e volte. Repare que ele pergunta se quer guardar, mas nesse caso guarda num geral, e pode guardar mais do que queremos. Por exemplo, podemos guardar a parte importante, e depois ficarmos a fazer algumas experiências ou contas, que não são para guardar...

Esta é a altura para aprender alguns comandos de limpeza e organização da casa. `ls` permite listar o que está no ambiente, e `rm` apagar.

```
> ls()
character(0)
> load("tudo.Rdata")
> ls()
[1] "bobagem"          "CondivSimples"  "onlyfb"         "SoFutebol"
```

Outra coisa que pode dar jeito, e que ocupa menos espaço, é manter uma história dos comandos, para a qual o R também tem um formato especial, `.Rhistory`. Pode dar-se um nome:

```
> savehistory("comandosTese.Rhistory")
```

Como deverá ser fácil de adivinhar, para os ir buscar numa futura sessão, basta invocar o comando dual `loadhistory`.

```
> loadhistory("comandosTese.Rhistory")
```

Em todos estes casos, é preciso reparar que, se nada diferente for dito, o R assume que estes objetos estão na diretoria de trabalho. Se não estiverem, tem de se incluir o nome completo do arquivo dentro das aspas, ou mudar para o sítio certo (com `setwd`) antes de o carregar.

5 Exercícios de estatística

(Estes exercícios são para ser feitos, explicados, e depois resolvidos, durante o próprio curso, por isso a informação aqui é mínima, propositadamente.)

5.1 Referência ao medo em português

Leia o ficheiro `http://folk.uio.no/dssantos/cursoR/medo.txt`, para uma folha de registo chamada `medo`. (A partir de agora, todos os nomes dos arquivos pressupõem o mesmo caminho, `http://folk.uio.no/dssantos/cursoR/`, que não será mais repetido.)

Familiarize-se com essa folha de registo, e explore-a graficamente.

Crie uma nova coluna com a informação número de palavras de medo por 10.000 palavras, e crie a figura correspondente.

Existe diferença entre os gêneros textuais? Use métodos estatísticos para responder a esta pergunta. Para saber mais sobre o medo em português e em inglês, veja [7].

5.2 Estudo das cores em português

Vá buscar o arquivo `CondivCompleto.txt`.

Investigue se existem diferenças entre o número de palavras de cor utilizadas por década, por variante do português, e por tema.

Obtenha um modelo dos dados, e observe se existem dependências entre as variáveis independentes.

Verifique, no corpo `Vercial`, o que pode dizer ou descobrir acerca da evolução das cores, ou da sua cor preferida, usando o material de `vercial.txt`.

5.3 Estudo de diferenças lexicais entre as variantes

Para que casos é possível demonstrar estatisticamente que as palavras são usadas de forma diferente nas duas variantes?

Exemplos reais, retirados de [11] e dos corpos do AC/DC, e presentes em `DifsLex.txt` e `DifsLex2.txt`.

Medida da incerteza e do intervalo de confiança.

5.4 Diferença entre tamanhos de frases entre originais e traduções

As frases nos textos traduzidos têm sempre um tamanho superior ao das frases originais? Ou depende da língua? Ou doutra coisa qualquer?

Os dados são os tamanhos das frases (em termos de número de palavras) no `CorTrad` jornalístico, em `tamanhosCorTradJorn.txt`.

Nota: Quando uma frase é traduzida por várias, conta-se a soma das várias. Quando mais do que uma frase é convertida para uma só, conta-se a soma das várias.

Experimente retirar, dos seus testes, os casos de frases não traduzidas, de forma a não contarem em cálculos de tamanho.

5.5 Estudo de diferenças sintáticas entre mulheres e homens

Leia o ficheiro `austen.txt`, com dados retirados de [14].

Investigue se há diferenças no uso de adjetivos entre mulheres e homens, em geral, e nos casos de conversas entre géneros, nos livros de Jane Austen. Para esta última questão, use `austenCompleto.txt`.

5.6 Estudo de diferenças semânticas entre autores

Dados os valores sobre cores obtidos com o `COMPARA` e apresentados em [16], podemos afirmar que há diferenças estatisticamente significativas entre os autores nessa matéria?

Os arquivos relevantes chamam-se `autCoresCOMPARA.txt`, com a distribuição de cores pro autor, e `autAdjCOMPARA.txt`, com a distribuição de adjetivos por autor.

5.7 O uso dos tempos em português

Passemos agora para métodos exploratórios. Existem contextos em que certos tempos são muito mais frequentes do que outros? Existem verbos que atraem certos tempos e rejeitam outros?

Estude o problema usando as listas de frequências de tempo e aspeto que se encontram no arquivo `tempos.txt`, inicialmente apresentadas em [12].

5.8 A invisibilidade do tradutor

Será que é possível detetar diferenças entre tradutores masculinos e femininos? E portanto determinar, com alguma segurança, qual o sexo do tradutor?

Dado um conjunto de características obtidas para uma série de textos traduzidos de várias línguas para inglês, [13] tentaram averiguar se era possível descobrir o sexo do tradutor. Tentemos repetir essa pesquisa, embora de uma forma simplificada. Os dados estão em `sexodotradutor.txt`.

6 Leituras recomendadas

Na secção de referências incluo vários livros escritos para apresentar o R a pessoas que trabalham com língua, do qual destaco o livro de Baayen [1] como o mais completo, assim como outros livros que são interessantes quer para aprender estatística (com R) quer para aprender a aplicação da estatística à linguística ou ao processamento computacional da língua.

Também incluo alguma documentação de referência do R, e materiais de apoio criados por outros, que se encontrem acessíveis na rede.

Faço também uma referência especial à tese do Stefan Evert [4], por ser extremamente clara e didática no que se refere ao estudo de co-ocorrências e de colocações.

O caminho de cada um através desta floresta de documentação, não o pretendo predeterminar. Bom passeio, e boa colheita... de ideias, métodos, e programas!

Referências

- [1] Harald Baayen. *Analyzing Linguistic Data: A practical introduction to Statistics using R*, Cambridge University Press, 2008.
- [2] R. H. Baayen. *languageR: Data sets and functions with "Analyzing Linguistic Data: A practical introduction to statistics*, 2011. <http://CRAN.R-project.org/package=languageR>. R package version 1.2.

- [3] Michael J. Crawley. *Statistics: An Introduction using R*, John Wiley and Sons, 2005.
- [4] Stefan Evert. *The Statistics of Word Cooccurrences: Word Pairs and Collocations*. Tese de doutoramento, Institut für maschinelle Sprachverarbeitung, University of Stuttgart, 2004.
- [5] Stefan Th. Gries. *Statistics for Linguistics with R: A Practical Introduction*, Mouton de Gruyter, 2009.
- [6] Keith Johnson. *Quantitative Methods in Linguistics*, Wiley-Blackwell, 2008.
- [7] Belinda Maia e Diana Santos. Who is afraid of... what? - In English and in Portuguese. Em Signe Oksefjell Ebeling, Jarle Ebeling e Hilde Hasselgård, editores, *Proceedings of ICAME 32*, 2012.
- [8] Michael P. Oakes. *Statistics for Corpus Linguistics*, Edinburgh University Press, 1998.
- [9] Michael P. Oakes e Meng Ji, editores. *Quantitative methods in corpus-based translation studies : a practical guide to descriptive translation research*, John Benjamims, 2012.
- [10] R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, 2010. <http://www.R-project.org>. ISBN 3-900051-07-0.
- [11] Tommaso Raso e Heliana Melo, editores. *C-oral-Brasil I Corpus de referência do português brasileiro falado informal*, Editora UFMG, 2012.
- [12] Diana Santos. The next step for the translation network. 2012. p. 35–52.
- [13] Miriam Schlesinger, Moshe Koppel, Noam Ordan e Brenda Malkiel. Markers of translator gender: Do they really matter? *Copenhagen Studies in Language*, 38:138–198, 2009.
- [14] Kari Anne Rand Schmidt. Male and female language in Jane Austen’s novels. Em Stig Johansson e Bjørn Tysdahl, editores, *Papers from the First Nordic Conference for English Studies (Oslo, 19-19 September, 1980)*, 1980. p. 198–210.
- [15] Tom Short. *R Reference Card*, 2004. <http://cran.r-project.org/doc/contrib/Short-refcard.pdf>.

- [16] Rosário Silva, Susana Inácio e Diana Santos. Colouring COMPARA: contrastive and monolingual colour studies in English and Portuguese. Apresentação na American Association for Corpus Linguistics (AACL) (Provo, USA, March 2008), 2008. http://www.linguateca.pt/Diana/download/Silvaetal_AACL2008_final.pdf.
- [17] Luís Torgo. *Acetatos usados nas aulas*, sem data. <http://www.liaad.up.pt/~ltorgo/Ensino/FEP/>. Faculdade de Economia, Universidade do Porto.