28/10/2025 18:20 1/5 Importazione dei testi

Importazione dei testi

Importazione di file di testo in R, in formati adatti all'uso con principali pacchetti per il text mining e l'analisi testuale.

I **dati degli esempi** sono disponibili a questo link. (in aggiornamento).

Vedi anche: Strumenti per l'analisi testuale e il text mining con R

Importazione di un solo testo

Con il pacchetto *readtext*, il file di testo viene importato come dataframe (funzione readtext()).

```
library(tidyverse)
library(readtext)
# sepolcri
sepolcri <- readtext("dati/sepolcri.txt")</pre>
str(sepolcri)
Classes 'readtext' and 'data.frame': 1 obs. of 2 variables:
 $ doc id: chr "sepolcri.txt"
 $ text : chr "All'ombra de' cipressi e dentro l'urne\nConfortate di
pianto è forse il
sonno\nDella morte men duro? Ove più il"| truncated
# articolo
articolo <- readtext("dati/ansa 2020-02-05.txt")</pre>
str(articolo)
Classes 'readtext' and 'data.frame':
                                        1 obs. of 2 variables:
$ doc_id: chr "ansa_2020-02-05.txt"
 $ text : chr "##TITOLO Sanremo: Fiorello al festival mai più
```

I nomi dei campi sono quelli riconosciuti per default dalle funzioni di Quanteda, Tidytext e tm, ovvero:

• doc id: identificativo di documento.

ha colpito Achille "| __truncated__

• text: campo con il testo da analizzare.

ospite, magari in gara\n\n##SOTTOTITOLO \"Mi

È possibile importare anche testi organizzati in dataset (formati: csv, tab, tsv) e documenti di Word e di OpenOffice. Per l'esattezza, i formati di documento accettati sono: txt, odt, doc, docx, ma anche rtf,pdf, json, html, e xml.

Importare più testi organizzati in una cartella

2 02 Alberto Urso.txt "\"Alberto Ur\"..."

Vengono così importati tutti i file della cartella, e nel campo doc_id viene conservato il nome dei file.

Le variabili relative ai singoli documenti possono essere estrapolate dal nome o dal percorso del file (le variabili saranno i nomi delle cartelle: nell'esempio che segue, docvarsfrom = "filenames").

Le variabili possono anche essere importate da un file csv (con un identificativo che associ le variabili al testo).

Il fatto che i testi siano organizzati in un dataframe permette di usare le funzioni per la gestione dei dati, dopo l'importazione. Ad esempio, potremmo voler trasformare una variabile "data" nel più appropriato formato *Date*:

28/10/2025 18:20 3/5 Importazione dei testi

```
articolo$data <- as.Date(articolo$data)
articolo</pre>
```

Importare un dataset con i testi

Nel caso di testi brevi o molto brevi, come le risposte aperte a un questionario, o i messaggi di Twitter, i dati sono generalmente raccolti e organizzati in un dataframe. In questo caso:

- non useremo readtext(), ma le funzioni del pacchetto readr (quindi potremo importare il dataset anche dal menu di RStudio);
- potremo strutturare le variabili anche prima dell'importazione.

Vediamo qui il dataset contenente le canzoni di Sanremo 2021, in formato tab delimited (tsv, anche se con estensione .csv):

```
sanremo21 <- read_tsv("dati/sanremo2021.csv")
head(sanremo21)</pre>
```

```
# A tibble: 6 x 3
  cantante
                     titolo
                                         testo
 <chr>
                     <chr>
                                         <chr>
1 Aiello
                     0ra
                                          "Ora ora ora\nMi parli
come allora~
2 Annalisa Scarrone
                     Dieci
                                          "Cos'è che ti ho
promesso\nNon so\nNon~
3 Arisa
                     Potevi Fare Di Più
                                         "Lasciarsi adesso non fa più
male non ~
                     E Invece Sì
                                          "Le metropolitane vanno
4 Bugo
molto veloci\n~
5 COLAPESCEDIMARTINO Musica Leggerissima "Se fosse un'orchestra a
parlare per n~
6 Coma Cose
                     Fiamme Negli Occhi
                                          "Quando ti sto vicino
sento\nChe a vol~
```

L'importazione può essere effettuata anche da RStudio (Import Dataset).

Organizziamo i nomi dei campi secondo gli standard visti sopra:

```
# A tibble: 6 x 4
  doc id
                      text
                                                          cantante
titolo
 <chr>
                      <chr>
                                                          <chr>
<chr>
                      "Ora ora ora ora\nMi parli come a~ Aiello
1 0ra
0ra
2 Dieci
                      "Cos'è che ti ho promesso\nNon so~ Annalisa S~
Dieci
                      "Lasciarsi adesso non fa più male~ Arisa
3 Potevi Fare Di Più
Potevi Far~
4 E Invece Sì
                      "Le metropolitane vanno molto vel~ Bugo
E Invece Sì
5 Musica Leggerissima "Se fosse un'orchestra a parlare ~ COLAPESCED~
Musica Leg~
6 Fiamme Negli Occhi
                      "Quando ti sto vicino sento\nChe ~ Coma Cose
Fiamme Neg~
```

Script di esempio

E' possibile scaricare ed eseguire lo script dell'esempio:

import_testi.R

```
library(tidyverse)
library(readtext)

# sepolcri
sepolcri <- readtext("dati/sepolcri.txt")

# articolo
articolo <- readtext("dati/ansa_2020-02-05.txt")
articolo$data <- as.Date(articolo$data)

# canzoni</pre>
```

28/10/2025 18:20 5/5 Importazione dei testi

Analisi testuale, Importazione, Gestione dei dati

From

https://www.agnesevardanega.eu/wiki/ - Ricerca Sociale con R

Permanent link:

https://www.agnesevardanega.eu/wiki/r/analisi-testuale/importazione testi

Last update: 13/08/2025 10:35

