

Package ‘enc’

October 13, 2022

Title Portable Tools for 'UTF-8' Character Data

Date 2019-12-19

Version 0.2.2

Description Implements an S3 class for storing 'UTF-8' strings, based on regular character vectors. Also contains routines to portably read and write 'UTF-8' encoded text files, to convert all strings in an object to 'UTF-8', and to create character vectors with various encodings.

Depends R (>= 3.1)

Imports methods

Suggests digest, pillar, readr, rlang, testthat, withr

License GPL-3

Encoding UTF-8

LazyData true

BugReports <https://github.com/krlmlr/enc/issues>

URL <https://github.com/krlmlr/enc>

RoxygenNote 7.0.2

NeedsCompilation yes

Author Kirill Müller [aut, cre]

Maintainer Kirill Müller <krlmlr+r@mailbox.org>

Repository CRAN

Date/Publication 2019-12-19 16:50:13 UTC

R topics documented:

enc-package	2
encoding	2
read_lines_enc	3
to_encoding	4
transform_lines_enc	5
utf8	6
write_lines_enc	7

Index**8**

enc-package *enc: Portable Tools for 'UTF-8' Character Data*

Description

Implements an S3 class for storing 'UTF-8' strings, based on regular character vectors. Also contains routines to portably read and write 'UTF-8' encoded text files, to convert all strings in an object to 'UTF-8', and to create character vectors with various encodings.

Author(s)

Maintainer: Kirill Müller <krlmlr+r@mailbox.org>

See Also

Useful links:

- <https://github.com/krlmlr/enc>
- Report bugs at <https://github.com/krlmlr/enc/issues>

encoding *Encoding information*

Description

The encoding function returns "ASCII" if the entire value consists of ASCII symbols only, and works identically to `Encoding()` otherwise. The `all_utf8` function is an efficient variant of `all(encoding(x) %in% c("ASCII", "UTF-8"))`.

Usage

```
encoding(x)
```

```
all_utf8(x)
```

Arguments

x A character vector.

Examples

```
encoding("a")
encoding("\u00fc")
all_utf8(enc2utf8(c("a", "\u00fc")))

# Platform-dependent:
all_utf8(enc2native(c("a", "\u00fc")))
```

read_lines_enc	<i>Reads from a text file</i>
----------------	-------------------------------

Description

This function is a drop-in replacement for [readLines\(\)](#) from disk files. It always returns text in the UTF-8 encoding and never warns on missing EOL on the last line.

[try_read_lines_enc\(\)](#) is a variant that returns an empty character vector on error, with a warning.

Usage

```
read_lines_enc(  
  path,  
  file_encoding = "UTF-8",  
  n = -1L,  
  ok = TRUE,  
  skipNul = FALSE  
)
```

```
try_read_lines_enc(  
  path,  
  file_encoding = "UTF-8",  
  n = -1L,  
  ok = TRUE,  
  skipNul = FALSE  
)
```

Arguments

path	Path to the file
file_encoding	The encoding to assume for the input file.
n	integer. The (maximal) number of lines to read. Negative values indicate that one should read up to the end of input on the connection.
ok	logical. Is it OK to reach the end of the connection before $n > 0$ lines are read? If not, an error will be generated.
skipNul	logical: should nuls be skipped?

See Also

[readr::read_lines\(\)](#) for a faster alternative.

Other file functions: [transform_lines_enc\(\)](#), [write_lines_enc\(\)](#)

to_encoding

Deep conversion to an encoding

Description

Converts all characters directly or indirectly contained in an object to a specific encoding. This works even if the encoding is different in the elements of a character vector.

Usage

```
to_encoding
to_utf8(x, ...)
to_native(x, ...)
to_latin1(x, ...)
to_alien(x, ...)

## S3 method for class 'character'
to_encoding(x, ..., converter)
```

Arguments

x	A character vector.
...	passed on to methods
converter	A function that accepts a character value as first argument and returns a (possibly classed) character with the desired encoding

Format

An object of class function of length 1.

Details

to_utf8 converts to UTF-8, using the `utf8()` class where possible. Implemented as `to_encoding(x, as_utf8)`

to_native converts to the native encoding. Implemented as `to_encoding(x, enc2native)` on Windows and as `to_encoding(x, as_utf8)` on Linux and OS X

to_latin1 converts to the latin-1 encoding

to_alien converts to the "other" encoding, i.e., UTF-8 on Windows and latin-1 on Linux and OS X.

See Also

- `rlang::as_utf8_character()` and `iconv()` for different ways to convert character vectors to Unicode

Examples

```
to_utf8(letters)
to_utf8(iris)
class(levels(to_utf8(iris)$Species))
```

transform_lines_enc *Transform a text file*

Description

Reads a file from disk, applies a function on the contents, and optionally writes the file back if different. The line ending separator of the input file is used if it can be read and contains at least one, otherwise `native_eol()` is used.

Usage

```
transform_lines_enc(
  path,
  fun,
  file_encoding = "UTF-8",
  ok = TRUE,
  skipNul = FALSE,
  write_back = TRUE,
  verbose = interactive()
)
```

Arguments

<code>path</code>	A vector of file paths.
<code>fun</code>	A function that returns a character vector.
<code>file_encoding</code>	The encoding to assume for the input file.
<code>ok</code>	logical. Is it OK to reach the end of the connection before $n > 0$ lines are read? If not, an error will be generated.
<code>skipNul</code>	logical: should nuls be skipped?
<code>write_back</code>	Should the results of the transformation be written back to the file?
<code>verbose</code>	Should the function show a message with a list of changed files?

Value

A named logical vector of the same length as `path` that indicates if a file has changed (TRUE or FALSE), or if an error occurred (NA)

See Also

Other file functions: [read_lines_enc\(\)](#), [write_lines_enc\(\)](#)

 utf8

A simple class for storing UTF-8 strings

Description

The values are stored as a [character\(\)](#) vector. On construction, the [enc2utf8\(\)](#) function is called on the input. Subsetting and concatenation operations on an object of this class return an object of this class again. Calls to [Encoding<-\(\)](#) are not intercepted.

Usage

```
utf8(x = character())

is_utf8(x)

as_utf8(x, ...)

## Default S3 method:
as_utf8(x, ...)

## S3 method for class ``NULL``
as_utf8(x, ...)

## S3 method for class 'character'
as_utf8(x, ...)

## S3 method for class 'utf8'
as_utf8(x, ...)

## S3 method for class 'utf8'
as.character(x, ...)

## S3 method for class 'utf8'
as.data.frame(
  x,
  row.names = NULL,
  optional = FALSE,
  ...,
  nm = paste(deparse(substitute(x)), width.cutoff = 500L), collapse = " ")
)

## S3 method for class 'utf8'
format(x, ...)
```

```
## S3 method for class 'utf8'
print(x, ...)
```

Arguments

x	A vector
...	Arguments passed on to further methods.
row.names	NULL or a character vector giving the row names for the data frame. Missing values are not allowed.
optional	logical. If TRUE, setting row names and converting column names (to syntactic names: see make.names) is optional. Note that all of R's base package <code>as.data.frame()</code> methods use <code>optional</code> only for column names treatment, basically with the meaning of <code>data.frame(*, check.names = !optional)</code> . See also the <code>make.names</code> argument of the <code>matrix</code> method.
nm	Name of column in new data frame

Examples

```
utf8(letters)
utf8("ä")
utf8(iconv("ä", to = "latin1"))
```

write_lines_enc	<i>Writes to a text file</i>
-----------------	------------------------------

Description

This function is a drop-in replacement for `writelnLines()` from `disk` files. It always expects text in the UTF-8 encoding, and by default writes in the UTF-8 encoding with Unix line separators.

Usage

```
write_lines_enc(text, path, file_encoding = "UTF-8", sep = "\n")
```

Arguments

text	A character vector
path	Path to the file
file_encoding	The encoding for the output file.
sep	character string. A string to be written to the connection after each line of text.

See Also

[readr::write_lines\(\)](#) for a faster alternative.

Other file functions: [read_lines_enc\(\)](#), [transform_lines_enc\(\)](#)

Index

- * **datasets**
 - to_encoding, 4
- * **file functions**
 - read_lines_enc, 3
 - transform_lines_enc, 5
 - write_lines_enc, 7
- all_utf8(encoding), 2
- as.character.utf8(utf8), 6
- as.data.frame.utf8(utf8), 6
- as_utf8(utf8), 6
- character(), 6
- data.frame, 7
- enc-package, 2
- enc2utf8(), 6
- encoding, 2
- Encoding(), 2
- format.utf8(utf8), 6
- iconv(), 5
- is_utf8(utf8), 6
- make.names, 7
- native_eol(), 5
- print.utf8(utf8), 6
- read_lines_enc, 3, 6, 7
- readLines(), 3
- readr::read_lines(), 3
- readr::write_lines(), 7
- rlang::as_utf8_character(), 5
- to_alien(to_encoding), 4
- to_encoding, 4
- to_latin1(to_encoding), 4
- to_native(to_encoding), 4
- to_utf8(to_encoding), 4
- transform_lines_enc, 3, 5, 7
- try_read_lines_enc(read_lines_enc), 3
- utf8, 6
- utf8(), 4
- write_lines_enc, 3, 6, 7
- writeLines(), 7