

Package ‘audubon’

October 12, 2022

Title Japanese Text Processing Tools

Version 0.3.0

Description A collection of Japanese text processing tools for filling Japanese iteration marks, Japanese character type conversions, segmentation by phrase, and text normalization which is based on rules for the 'Sudachi' morphological analyzer and the 'NEologd' (Neologism dictionary for 'MeCab'). These features are specific to Japanese and are not implemented in 'ICU' (International Components for Unicode).

License Apache License (≥ 2)

URL <https://github.com/paithiov909/audubon>,
<https://paithiov909.github.io/audubon/>

BugReports <https://github.com/paithiov909/audubon/issues>

Depends R (≥ 2.10)

Imports dplyr, magrittr, memoise, purrr, readr, rlang ($\geq 0.4.11$),
stringi, tidyselect, V8

Suggests roxygen2, spelling, testthat ($\geq 3.0.0$)

Config/testthat/edition 3

Encoding UTF-8

Language en-US

LazyData true

RoxygenNote 7.2.1

NeedsCompilation no

Author Akiru Kato [cre, aut],
Koki Takahashi [cph] (Author of japanese.js),
Shuhei Iitsuka [cph] (Author of budoux),
Taku Kudo [cph] (Author of TinySegmenter)

Maintainer Akiru Kato <paithiov909@gmail.com>

Repository CRAN

Date/Publication 2022-07-22 17:20:02 UTC

R topics documented:

| | |
|---------------------|-----------|
| get_dict_features | 2 |
| ngram_tokenizer | 3 |
| pack | 3 |
| polano | 4 |
| prettify | 5 |
| read_rewrite_def | 5 |
| strj_fill_iter_mark | 6 |
| strj_hiraganize | 6 |
| strj_katakanize | 7 |
| strj_normalize | 8 |
| strj_rewrite_as_def | 8 |
| strj_romanize | 9 |
| strj_segment | 10 |
| strj_tinyseg | 11 |
| strj_tokenize | 12 |
| strj_transcribe_num | 13 |
| Index | 14 |

| | |
|-------------------|----------------------------------|
| get_dict_features | <i>Get dictionary's features</i> |
|-------------------|----------------------------------|

Description

Returns dictionary's features. Currently supports "unic17" (2.1.2 src schema), "unic26" (2.1.2 bin schema), "unic29" (schema used in 2.2.0, 2.3.0), "cc-cedict", "ko-dic" (mecab-ko-dic), "naist11", "sudachi", and "ipa".

Usage

```
get_dict_features(
  dict = c("ipa", "unic17", "unic26", "unic29", "cc-cedict", "ko-dic", "naist11",
    "sudachi")
)
```

Arguments

| | |
|------|---|
| dict | Character scalar; one of "ipa", "unic17", "unic26", "unic29", "cc-cedict", "ko-dic", "naist11", or "sudachi". |
|------|---|

Value

A character vector.

See Also

See also ['CC-CEDICT-MeCab'](#), and ['mecab-ko-dic'](#).

Examples

```
get_dict_features("ipa")
```

| | |
|-----------------|-------------------------|
| ngram_tokenizer | <i>Ngrams tokenizer</i> |
|-----------------|-------------------------|

Description

Make an ngram tokenizer function.

Usage

```
ngram_tokenizer(n = 1L)
```

Arguments

| | |
|---|----------|
| n | Integer. |
|---|----------|

Value

ngram tokenizer function

| | |
|------|---|
| pack | <i>Pack prettified data.frame of tokens</i> |
|------|---|

Description

Packs a prettified data.frame of tokens into a new data.frame of corpus, which is compatible with the Text Interchange Formats.

Usage

```
pack(tbl, pull = "token", n = 1L, sep = "-", .collapse = " ")
```

Arguments

| | |
|-----------|--|
| tbl | A prettified data.frame of tokens. |
| pull | Column to be packed into text or ngrams body. Default value is token. |
| n | Integer internally passed to ngrams tokenizer function created of audubon::ngram_tokenizer() |
| sep | Character scalar internally used as the concatenator of ngrams. |
| .collapse | This argument is passed to stringi::stri_join(). |

Value

A data.frame.

Text Interchange Formats (TIF)

The Text Interchange Formats (TIF) is a set of standards that allows R text analysis packages to target defined inputs and outputs for corpora, tokens, and document-term matrices.

Valid data.frame of tokens

The prettified data.frame of tokens here is a data.frame object compatible with the TIF.

A TIF valid data.frame of tokens are expected to have one unique key column (named doc_id) of each text and several feature columns of each tokens. The feature columns must contain at least token itself.

See Also

<https://github.com/ropensci/tif>

Examples

```
pack(strj_tokenize(polano[1:5], format = "data.frame"))
```

| | |
|--------|---|
| polano | <i>Whole text of 'Porano no Hiroba' written by Miyazawa Kenji from Aozora Bunko</i> |
|--------|---|

Description

Whole text of 'Porano no Hiroba' written by Miyazawa Kenji from Aozora Bunko

Usage

```
polano
```

Format

An object of class character of length 899.

Details

A dataset containing the text of Miyazawa Kenji's novel "Porano no Hiroba" which was published in 1934, the year after Kenji's death. Copyright of this work has expired since more than 70 years have passed after the author's death.

The UTF-8 plain text is sourced from <https://www.aozora.gr.jp/cards/000081/card1935.html> and is cleaned of meta data.

Source

https://www.aozora.gr.jp/cards/000081/files/1935_ruby_19924.zip

Examples

```
head(polano)
```

```
prettify          Prettify tokenized output
```

Description

Prettify tokenized output

Usage

```
prettify(df, into = get_dict_features("ipa"), col_select = seq_along(into))
```

Arguments

`df` A data.frame that has feature column to be prettified.
`into` Character vector that is used as column names of features.
`col_select` Character or integer vector that will be kept in prettified result.

Value

A data.frame.

```
read_rewrite_def  Read a rewrite.def file
```

Description

Read a rewrite.def file

Usage

```
read_rewrite_def(
  def_path = system.file("def/rewrite.def", package = "audubon")
)
```

Arguments

`def_path` Character scalar; path to the rewriting definition file.

Value

A list.

Examples

```
str(read_rewrite_def())
```

strj_fill_iter_mark *Fill Japanese iteration marks*

Description

Fills Japanese iteration marks (Odori-ji) with their previous characters if the element has more than 5 characters.

Usage

```
strj_fill_iter_mark(text)
```

Arguments

text Character vector.

Value

A character vector.

Examples

```
strj_fill_iter_mark(c(
  "\u3042\u3044\u3046\u309d\u3003\u304b\u304d",
  "\u91d1\u5b50\u307f\u3059\u309e",
  "\u306e\u305f\u308a\u3033\u3035\u304b\u306a",
  "\u3057\u308d\u2033\u203c\u3068\u3057\u305f"
))
```

strj_hiraganize *Hiraganize Japanese characters*

Description

Converts Japanese katakana to hiragana. It is almost similar to `stringi::stri_trans_general(text, "kana-hira")`, however, this implementation can also handle some additional symbols such as Japanese kana ligature (aka. goryaku-gana).

Usage

```
strj_hiraganize(text)
```

Arguments

text Character vector.

Value

A character vector.

Examples

```
strj_hiraganize(
  c(
    paste0(
      "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
      "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
      "\u3068\u304a\u3063\u305f\u98a8"
    ),
    "\u677f\u57a3\u6b7b\u30b9\u0002a708"
  )
)
```

| | |
|-----------------|---------------------------------------|
| strj_katakanize | <i>Katakanize Japanese characters</i> |
|-----------------|---------------------------------------|

Description

Converts Japanese hiragana to katakana. It is almost similar to `stringi::stri_trans_general(text, "hira-kana")`, however, this implementation can also handle some additional symbols such as Japanese kana ligature (aka. goryaku-gana).

Usage

```
strj_katakanize(text)
```

Arguments

`text` Character vector.

Value

A character vector.

Examples

```
strj_katakanize(
  c(
    paste0(
      "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
      "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
      "\u3068\u304a\u3063\u305f\u98a8"
    ),
    "\u672c\u65e5\u309f\u304b\u304d\u6c37\u89e3\u7981"
  )
)
```

strj_normalize *Convert text following the rules of 'NEologd'*

Description

Converts characters into normalized style following the rule that is recommended by the Neologism dictionary for 'MeCab'.

Usage

```
strj_normalize(text)
```

Arguments

text Character vector to be normalized.

Value

A character vector.

See Also

<https://github.com/neologd/mecab-ipadic-neologd/wiki/Regexp.ja>

Examples

```
strj_normalize(  
  paste0(  
    "\u2015\u2015\u5357\u30a2\u30eb\u30d7\u30b9",  
    "\u306e\u3000\u5929\u7136\u6c34-\u3000\u3033",  
    "\uff50\u41\u52\u4b\u49\u4e\u47*",  
    "\u3000\u2c\u45\u4d\u4f\u4e+",  
    "\u3000\u30ec\u30e2\u30f3\u4e00\u7d5e\u308a"  
  )  
)
```

strj_rewrite_as_def *Rewrite text using rewrite.def*

Description

Rewrite text using rewrite.def

Usage

```
strj_rewrite_as_def(text, as = read_rewrite_def())
```


Arguments

text Character vector to be normalized.
as List.

Value

A character vector.

Examples

```
strj_rewrite_as_def(
  paste0(
    "\u2015\u2015\u5357\u30a2\u30eb",
    "\u30d7\u30b9\u306e\u3000\u5929",
    "\u7136\u6c34-\u3000\u3033\u3035",
    "\uff41\u3035\u303b\u303d\u303e\u303f*",
    "\u3000\u303c\u303e\u303f\u3040\u3041\u3042",
    "\u3000\u303c\u303e\u303f\u3040\u3041\u3042"
  )
)
strj_rewrite_as_def(
  "\u60e1\u3068\u5047\u9762\u306e\u30eb\u30fc\u30eb",
  read_rewrite_def(system.file("def/kyuji.def", package = "audubon"))
)
```

strj_romanize

Romanize Japanese Hiragana and Katakana

Description

Romanize Japanese Hiragana and Katakana

Usage

```
strj_romanize(
  text,
  config = c("wikipedia", "traditional hepburn", "modified hepburn", "kunrei", "nihon")
)
```

Arguments

text Character vector. If elements are composed of except but hiragana and katakana letters, those letters are dropped from the return value.
config Configuration used to romanize. Default is wikipedia.

Details

There are several ways to romanize Japanese. Using this implementation, you can convert hiragana and katakana as 5 different styles; the wikipedia style, the traditional hepburn style, the modified hepburn style, the kunrei style, and the nihon style.

Note that all of these styles return a slightly different form of `stringi::stri_trans_general(text, "Any-latn")`.

Value

A character vector.

See Also

<https://github.com/hakatashi/japanese.js#japaneseromanizetext-config>

Examples

```
strj_romanize(
  paste0(
    "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
    "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
    "\u3068\u304a\u3063\u305f\u98a8"
  )
)
```

strj_segment

Segment text into phrases

Description

Segments Japanese text into several phrases using the 'google/budoux' JavaScript module or 'Tiny-Segmenter.js'.

Usage

```
strj_segment(
  text,
  format = c("list", "data.frame"),
  engine = c("budoux", "tinyseg"),
  split = FALSE
)
```

Arguments

| | |
|--------|--|
| text | Character vector to be tokenized. |
| format | Output format. Choose list or data.frame. |
| engine | Splitting engine to be used. |
| split | Logical. If true, the function splits the vector into some sentences using <code>stringi::stri_split_boundar</code> = "sentence") before tokenizing. |

Value

List or data.frame.

Examples

```
strj_segment(
  paste0(
    "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
    "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
    "\u3068\u304a\u3063\u305f\u98a8"
  )
)
strj_segment(
  paste0(
    "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
    "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
    "\u3068\u304a\u3063\u305f\u98a8"
  ),
  format = "data.frame"
)
```

strj_tinyseg

Segment text into phrases

Description

An alias of `strj_segment(engine = "tinyseg")`.

Usage

```
strj_tinyseg(text, format = c("list", "data.frame"), split = FALSE)
```

Arguments

| | |
|---------------------|---|
| <code>text</code> | Character vector to be tokenized. |
| <code>format</code> | Output format. Choose <code>list</code> or <code>data.frame</code> . |
| <code>split</code> | Logical. If true, the function splits vectors into some sentences using <code>stringi::stri_split_boundaries = "sentence"</code>) before tokenizing. |

Value

A list or data.frame.

Examples

```

strj_tinyseg(
  paste0(
    "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
    "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
    "\u3068\u304a\u3063\u305f\u98a8"
  )
)
strj_tinyseg(
  paste0(
    "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
    "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
    "\u3068\u304a\u3063\u305f\u98a8"
  ),
  format = "data.frame"
)

```

strj_tokenize

Split text into tokens

Description

Splits text into several tokens using specified tokenizer.

Usage

```

strj_tokenize(
  text,
  format = c("list", "data.frame"),
  engine = c("stringi", "mecab", "sudachipy"),
  rcpath = NULL,
  mode = c("C", "B", "A"),
  split = FALSE
)

```

Arguments

| | |
|--------|---|
| text | Character vector to be tokenized. |
| format | Output format. Choose list or data.frame. |
| engine | Tokenizer name. Choose one of 'stringi', 'mecab', or 'sudachipy'. Note that the specified tokenizer is installed and available when you use 'mecab' or 'sudachipy'. |
| rcpath | Path to a setting file for 'MeCab' or 'sudachipy' if any. |
| mode | Splitting mode for 'sudachipy'. |
| split | Logical. If true, the function splits the vector into some sentences using stringi::stri_split_boundar = "sentence") before tokenizing. |

Value

A list or data.frame.

Examples

```
strj_tokenize(
  paste0(
    "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
    "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
    "\u3068\u304a\u3063\u305f\u98a8"
  )
)
strj_tokenize(
  paste0(
    "\u3042\u306e\u30a4\u30fc\u30cf\u30c8",
    "\u30fc\u30f4\u30a9\u306e\u3059\u304d",
    "\u3068\u304a\u3063\u305f\u98a8"
  ),
  format = "data.frame"
)
```

strj_transcribe_num *Transcribe Arabic to Kansuji*

Description

Transcribes Arabic integers to Kansuji with auxiliary numerals.

Usage

```
strj_transcribe_num(int)
```

Arguments

int Integers.

Details

As its implementation is limited, this function can only transcribe numbers up to trillions. In case you convert much bigger numbers, try to use the 'arabic2kansuji' package.

Value

A character vector.

Examples

```
strj_transcribe_num(c(10L, 31415L))
```

Index

* datasets

polano, 4

get_dict_features, 2

ngram_tokenizer, 3

pack, 3

polano, 4

prettify, 5

read_rewrite_def, 5

strj_fill_iter_mark, 6

strj_hiraganize, 6

strj_katakanize, 7

strj_normalize, 8

strj_rewrite_as_def, 8

strj_romanize, 9

strj_segment, 10

strj_tinyseg, 11

strj_tokenize, 12

strj_transcribe_num, 13