

# Package: RmecabKo (via r-universe)

November 11, 2024

**Type** Package

**Title** Korean User Interface for MeCab in R

**Version** 0.1.7.0

**Author** Junhewk Kim

**Maintainer** Junhewk Kim <junhewk.kim@gmail.com>

**Description** This package provides useful functions for text mining in Korean. It depends major POS analysis on 'RcppMeCab' package.

**Imports** Rcpp, RcppMeCab, stringr

**LinkingTo** Rcpp, RcppMeCab

**License** GPL

**RoxygenNote** 6.1.1

**Encoding** UTF-8

**LazyData** true

**Config/pak/sysreqs** make libicu-dev libmecab-dev

**Repository** <https://junhewk.r-universe.dev>

**RemoteUrl** <https://github.com/junhewk/rmecabko>

**RemoteRef** HEAD

**RemoteSha** ca5b22d2079565084dc05aa8b3c39b15f0994309

## Contents

install_dic . . . . .	2
install_mecab . . . . .	2
nouns . . . . .	3
pos . . . . .	4
RmecabKo . . . . .	5
token_morph . . . . .	6
token_ngrams . . . . .	7
words . . . . .	8

<b>Index</b>	<b>9</b>
--------------	----------

---

install_dic	<i>Install Mecab-Ko-Dic in Linux and Mac OSX.</i>
-------------	---

---

**Description**

install\_dic installs Mecab-Ko-Dic.

**Usage**

```
install_dic()
```

**Details**

This code checks and installs Mecab-Ko-Dic in Linux and Mac OSX. This is essential for using custom-defined user dictionary. Installing Mecab-Ko-Dic needs system privileges, because it uses 'make install' to build from source and install it to system.

**Value**

None. The function will halt when the current operation system is not Linux or Mac OSX, or Mecab-Ko-Dic is installed already.

See examples in [Github](#).

**Examples**

```
## Not run:  
install_dic()  
  
## End(Not run)
```

---

install_mecab	<i>Install mecab-ko-msvc and mecab-ko-dic-msvc</i>
---------------	--

---

**Description**

install\_mecab installs Mecab-Ko-MSVC and Mecab-Ko-Dic-MSVC.

**Usage**

```
install_mecab(mecabLocation)
```

**Arguments**

mecabLocation a directory to install Mecab-Ko-MSVC and Mecab-Ko-Dic-MSVC.

**Details**

This code checks and installs Mecab-Ko-MSVC and Mecab-Ko-Dic-MSVC in user specified directory. Windows only.

**Value**

None. The function will halt when the current operation system is not Windows, or /mecabLocation/mecab.exe exists.

See examples in [Github](#).

**Examples**

```
## Not run:  
install_mecab("D:/Rlibs/mecab")  
  
## End(Not run)
```

---

nouns

*Noun extractor by mecab-ko*

---

**Description**

nouns returns nouns extracted from Korean phrases.

**Usage**

```
nouns(sentence, sys_dic = "", user_dic = "", parallel = FALSE)
```

**Arguments**

phrase            A character vector or character vectors.

**Details**

Noun extraction is used for many Korean text analysis algorithms. The function coerces input to UTF-8.

**Value**

List of nouns will be returned. Element name of the list are original phrases.

See examples in [Github](#).

**Examples**

```
## Not run:
nouns(c("Some Korean Phrases"))

## End(Not run)
```

---

 pos

*POS tagging by mecab-ko*


---

**Description**

pos returns part-of-speech (POS) tagged morpheme of Korean phrases.

**Usage**

```
pos(sentence, join = TRUE, format = c("list", "data.frame"),
     sys_dic = "", user_dic = "", parallel = FALSE)
```

**Arguments**

sentence	Character vector.
join	Boolean to determine providing POS tags with the morphemes or not. The default value is TRUE.
format	A data type for the result. The default value is "list". You can set this to "data.frame" to get a result as data frame format.
sys_dic	A location of system MeCab dictionary. The default value is "".
user_dic	A location of user-specific MeCab dictionary. The default value is "".
parallel	Boolean to determine using parallel analyzing. The default value is FALSE.

**Details**

This is a basic function of part-of-speech tagging by mecab-ko. The function coerces input to UTF-8.

**Value**

List of POS tagged morpheme will be returned in conjoined character vector form. Element name of the list are original phrases. If join=FALSE, it returns list of morpheme with named with tags.

See examples in [Github](#).

**Examples**

```
## Not run:
pos(c("Some Korean Phrases"))
pos(c("Some Korean Phrases"), join=FALSE)

## End(Not run)
```

RmecabKo

*Rcpp Wrapper for Eunjeon Project***Description**

The mecab-ko and mecab-ko-dic is based on a C++ library, and POS tagging with them is useful when the spacing of source text is not correct. For integrating mecab-ko with R, Rcpp package is used for providing the basic framework.

**Details**

It is based on the Eunjeon Project. For Mac OSX and Linux, You need to install mecab-ko and mecab-ko-dic before install this package in R. mecab-ko: <https://bitbucket.org/eunjeon/mecab-ko> mecab-ko-dic: <https://bitbucket.org/eunjeon/mecab-ko-dic> In Windows, `install_mecab(mecabLocation)` function will install mecab-ko-msvc and mecab-ko-dic-msvc in user specified directory. It is operated by system command and file I/O, the speed of the analysis is slow compared to the Linux-based operating system.

**Author(s)**

Junhewk Kim

**References**

- [Eunjeon project](#)
- [Wonsup Yoon, mecab-ko VC++ builds at https://github.com/Pusnow/mecab-ko-msvc, https://github.com/Pusnow/mecab-ko-dic-msvc](#)

**Examples**

```
## Not run:
# install.packages("devtools")
devtools::install_github("junhewk/RmecabKo")
# On Windows platform only
install_mecab("D:/Rlibs/mecab")

phrase <- # Some Korean character vectors

# For full POS tagging
pos(phrase)
```

```
# For noun extraction only
nouns(phrase)
# For tokenizing of selective morphemes
tokens_words(phrase)
# For n-grams tokenizing
tokens_ngram(phrase)

## End(Not run)
```

---

token\_morph

*Morpheme tokenizer based on mecab-ko*


---

## Description

These tokenizer functions perform tokenization into full or selected morphemes, nouns.

## Usage

```
token_morph(phrase, strip_punct = FALSE, strip_numeric = FALSE)
```

```
token_words(phrase, strip_punct = FALSE, strip_numeric = FALSE)
```

```
token_nouns(phrase, strip_punct = FALSE, strip_numeric = FALSE)
```

## Arguments

phrase	A character vector or a list of character vectors to be tokenized into morphemes. If phrase is a character vector, it can be of any length, and each element will be tokenized separately. If phrase is a list of character vectors, each element of the list should be a one-item vector.
strip_punct	Bool. If you want to remove punctuations in the phrase, set this as TRUE.
strip_numeric	Bool. If you want to remove numbers in the phrase, set this as TRUE.

## Value

A list of character vectors containing the tokens, with one element in the list.

See examples in [Github](#).

## Examples

```
## Not run:
txt <- # Some Korean sentence

token_morph(txt)
token_words(txt, strip_punct = FALSE)
token_nouns(txt, strip_numeric = TRUE)
```

```
## End(Not run)
```

---

token_ngrams	<i>N-gram tokenizer based on mecab-ko</i>
--------------	---

---

## Description

This function tokenizes inputs into n-grams. For the developmental purpose, this function offers basic n-gram (or shingle n-gram) only. Other n-gram functionality will be added later. Punctuations and numerics are stripped for this tokenizer, because in Korean n-grams those are usually useless. N-gram function is based on the selective morpheme tokenizer (`token_words`), but you can select other tokenizer as well.

## Usage

```
token_ngrams(phrase, n = 3L, div = c("morph", "words", "nouns"),
  stopwords = character(), ngram_delim = " ")
```

## Arguments

phrase	A character vector or a list of character vectors to be tokenized into morphemes. If phrase is a character vector, it can be of any length, and each element will be tokenized separately. If phrase is a list of character vectors, each element of the list should be a one-item vector.
n	The number of words in the n-gram. This must be an integer greater than or equal to 1.
div	The token generator definition. The options are "morph", "words", and "nouns".
stopwords	Stopwords set to exclude tokens.
ngram_delim	The separator between words in an n-gram.

## Value

A list of character vectors containing the tokens, with one element in the list.

See examples in [Github](#).

## Examples

```
## Not run:
txt <- # Some Korean sentence

token_ngrams(txt)
token_ngrams(txt, n = 2)

## End(Not run)
```

---

words

*Words extractor by mecab-ko*

---

### **Description**

words returns full morphemes extracted from Korean phrases.

### **Usage**

```
words(phrase)
```

### **Arguments**

phrase            Character vector.

### **Details**

It is based on Mecab-Ko POS classification. Full morphemes are consisted with The function coerces input to UTF-8.

### **Value**

List of full morphemes will be returned.

See examples in [Github](#).

### **Examples**

```
## Not run:  
words(c("Some Korean Phrases"))
```

```
## End(Not run)
```



# Index

- \* **Korean**

- RmecabKo, [5](#)

- \* **nlp**

- RmecabKo, [5](#)

- \* **tagger**

- RmecabKo, [5](#)

[install\\_dic](#), [2](#)

[install\\_mecab](#), [2](#)

[nouns](#), [3](#)

[pos](#), [4](#)

[RmecabKo](#), [5](#)

[RmecabKo-package \(RmecabKo\)](#), [5](#)

[token\\_morph](#), [6](#)

[token\\_ngrams](#), [7](#)

[token\\_nouns \(token\\_morph\)](#), [6](#)

[token\\_words \(token\\_morph\)](#), [6](#)

[words](#), [8](#)