

# Package ‘KoNLP’

July 2, 2014

**Maintainer** Heewon Jeon <madjakarta@gmail.com>

**License** GPL-3

**Title** Korean NLP Package

**Author** Heewon Jeon

**Description** Korean language processing package. Morphological analyzer, POS tagger, Keystroke converter, Hangul automata, Concordance, Mutual Information..

**SystemRequirements** Java (>= 1.6)

**URL** <https://github.com/haven-jeon/KoNLP>

**Version** 0.76.9

**Repository** CRAN

**Date/Publication** 2013-05-31 14:23:46

**Date** 2011-10-12

**Encoding** UTF-8

**Depends** R (>= 2.15.0), rJava (>= 0.9-0), utils (>= 2.14.0), stringr (>= 0.6.2), hash (>= 2.2.6), tau (>= 0.0-15), Sejong (>= 0.01)

**Collate** 'onLoad.R' 'manageDic.R' 'hangulUtils.R' 'koAnalyzerRun.R' 'tagdata.R' 'Concordances.R'

**NeedsCompilation** no

## R topics documented:

backupUsrDic . . . . .	2
concordance_file . . . . .	3
concordance_str . . . . .	4
convertHangulStringToJamos . . . . .	4
convertHangulStringToKeyStrokes . . . . .	5
convertTag . . . . .	5

editweights	6
extractNoun	6
HangulAutomata	7
is.ascii	7
is.hangul	8
is.jaeum	8
is.jamo	9
is.moeum	9
KtoS	10
mergeUserDic	10
MorphAnalyzer	11
mutualinformation	11
reloadAllDic	12
restoreUsrDic	12
SimplePos09	13
SimplePos22	14
statDic	14
StoK	15
tags	15
useSejongDic	15
useSystemDic	16

## Index 17

---

backupUsrDic	<i>use for backup current dic_user.txt</i>
--------------	--------------------------------------------

---

### Description

Utility function for backup dic\_user.txt file to backup directory.

### Usage

```
backupUsrDic(ask = TRUE)
```

### Arguments

ask	ask to confirm backup
-----	-----------------------

### Examples

```
## Not run:
## This codes can not be run
## if you don't have encoding system which can en/decode
## Hangul(ex) CP949, EUC-KR, UTF-8).
dicpath <- file.path(system.file(package="Sejong"), "dics", "handic.zip")
conn <- unz(dicpath, file.path("data","kE","dic_user2.txt"))
newdic <- read.table(conn, sep="\t", header=FALSE, fileEncoding="UTF-8", stringsAsFactors=FALSE)
mergeUserDic(newdic)
```

```
## backup merged new dictionary
backupUsrDic(ask=FALSE)
## restore from backup directory
restoreUsrDic(ask=FALSE)
## reloading new dictionary
reloadAllDic()
## End(Not run)
```

---

concordance\_file      *concordance for input text file*

---

### **Description**

returns concordance text for input file.

### **Usage**

```
concordance_file(filename, pattern,
  encoding = getOption("encoding"), span = 5)
```

### **Arguments**

filename	file name
pattern	patterns of central words
span	how many character will be produced around input pattern
encoding	filename's encoding

### **Author(s)**

Heewon Jeon

### **References**

Church, K. W. and Mercer, R. L. (1993). Introduction to the special issue on computational linguistics using large corpora. *Computational Linguistics*, 19(1):1-24.

concordance\_str      *concordance for input text vector*

---

**Description**

returns concordance text for input pattern and span.

**Usage**

```
concordance_str(string, pattern, span = 5)
```

**Arguments**

string	input text as character vector or single character
pattern	patterns of central words
span	how many character will be produced around input pattern

**Author(s)**

Heewon Jeon

**References**

Church, K. W. and Mercer, R. L. (1993). Introduction to the special issue on computational linguistics using large corpora. *Computational Linguistics*, 19(1):1-24.

---

convertHangulStringToJamos  
*conversion function Hangul string to Jamos*

---

**Description**

convert Hangul sentence to Jamos. Example will be shown in [github wiki](#).

**Usage**

```
convertHangulStringToJamos(hangul)
```

**Arguments**

hangul	Hangul string
--------	---------------

**Value**

Jamo sequences

---

convertHangulStringToKeyStrokes  
*conversion function Hangul string to keyStrokes*

---

### Description

Function can convert Hangul string to Keystrokes. Example will be shown in [github wiki](#).

### Usage

```
convertHangulStringToKeyStrokes(hangul,  
    isFullwidth = TRUE)
```

### Arguments

hangul	Hangul sentence
isFullwidth	specify returned character will be Fullwidth ASCII or Halfwidth ASCII

### Value

Keystroke sequence

---

convertTag            *tag name converter*

---

### Description

only support tag conversion between KAIST and Sejong tag set.

### Usage

```
convertTag(fromTag, toTag, tag)
```

### Arguments

fromTag	tag set name to convert from
toTag	desired tag set name
tag	tag name to search

---

editweights	<i>Keystroke misspell cost table</i>
-------------	--------------------------------------

---

**Description**

Keystroke misspell cost table

**Author(s)**

Heewon Jeon <madjakarta@gmail..com>

---

extractNoun	<i>Noun extractor for Hangul</i>
-------------	----------------------------------

---

**Description**

extract Nouns from Korean sentence uses Hannanum analyzer. see detail in [Hannanum](#). Example will be shown in [github wiki](#).

**Usage**

```
extractNoun(sentence)
```

**Arguments**

sentence	input
----------	-------

**Value**

Noun of sentence

**References**

Sangwon Park et al(2010). A Plug-In Component-based Korean Morphological Analyzer

---

HangulAutomata	<i>do Hangul automata</i>
----------------	---------------------------

---

**Description**

function to be used for converting to complete Hangul syllables from Jamo or Keystrokes. Example will be shown in [github wiki](#).

**Usage**

```
HangulAutomata(input, isKeystroke = F, isForceConv = F)
```

**Arguments**

input	to be processed mostly Jamo sequences
isKeystroke	boolean parameter to check input is keystroke or Jamo sequences
isForceConv	boolean parameter to force converting if input is not valid Jamo or keystroke sequences.

**Value**

complete Hangul syllable

---

is.ascii	<i>check if sentence is all ASCII</i>
----------	---------------------------------------

---

**Description**

Function checks with each character is ASCII

**Usage**

```
is.ascii(sentence)
```

**Arguments**

sentence	input characters
----------	------------------

**Value**

TRUE or FALSE

---

is.hangul	<i>check if sentence is all Hangul</i>
-----------	----------------------------------------

---

**Description**

Function checks if each character is Hangul or Jamo. Example will be shown in [github wiki](#).

**Usage**

```
is.hangul(sentence)
```

**Arguments**

sentence	input characters
----------	------------------

**Value**

TRUE or FALSE

---

is.jaeum	<i>check if sentence is all Jaeum</i>
----------	---------------------------------------

---

**Description**

Function checks with each character is Jaeum

**Usage**

```
is.jaeum(sentence)
```

**Arguments**

sentence	input characters
----------	------------------

**Value**

TRUE or FALSE



---

is.jamo

*check if sentence is all Jamo*

---

### **Description**

Function checks with each character is Jamo. Example will be shown in [github wiki](#).

### **Usage**

```
is.jamo(sentence)
```

### **Arguments**

sentence          input characters

### **Value**

TRUE or FALSE

---

is.moeum

*check if sentence is all Moeum*

---

### **Description**

Function checks with each character is Moeum

### **Usage**

```
is.moeum(sentence)
```

### **Arguments**

sentence          input characters

### **Value**

TRUE or FALSE

---

KtoS	<i>KAIST tag to Sejong tag</i>
------	--------------------------------

---

**Description**

KAIST tag to Sejong tag

**Author(s)**

Heewon Jeon <madjakarta@gmail..com>

---

mergeUserDic	<i>appending or replacing with new data.frame</i>
--------------	---------------------------------------------------

---

**Description**

appending new dictionary to current dictionary. replacing current dictionary with new dictionary.

**Usage**

```
mergeUserDic(newUserDic, append = TRUE, verbose = FALSE,
             ask = FALSE)
```

**Arguments**

newUserDic	new user dictionary as data.frame
append	append or replacing
verbose	see detail error logs
ask	ask to backup

**Examples**

```
## Not run:
## This codes can not be run
## if you don't have encoding system which can en/decode
## Hangeul(ex) CP949, EUC-KR, UTF-8).
dicpath <- file.path(system.file(package="Sejong"), "dics", "handic.zip")
conn <- unz(dicpath, file.path("data", "kE", "dic_user2.txt"))
newdic <- read.table(conn, sep="\t", header=FALSE, fileEncoding="UTF-8", stringsAsFactors=FALSE)
mergeUserDic(newdic)
## backup merged new dictionary
backupUsrDic(ask=FALSE)
## restore from backup directory
restoreUsrDic(ask=FALSE)
## reloading new dictionary
reloadAllDic()
## End(Not run)
```

---

MorphAnalyzer	<i>Hannanum morphological analyzer interface function</i>
---------------	-----------------------------------------------------------

---

**Description**

Do the morphological analysis, not doing pos tagging uses Hannanum analyzer. see details in [Hannanum](#). Example will be shown in [github wiki](#).

**Usage**

```
MorphAnalyzer(sentence)
```

**Arguments**

sentence	input
----------	-------

**Value**

result of analysis

**References**

Sangwon Park et al(2010). A Plug-In Component-based Korean Morphological Analyzer

---

mutualinformation	<i>mutual information for input text</i>
-------------------	------------------------------------------

---

**Description**

returns mutual information or t-scores for input text

**Usage**

```
mutualinformation(text, query = "",
  method = c("mutual", "tscores"))
```

**Arguments**

text	input character vector
method	for calculations('mutual' or 't-scores')
query	term to get information

**Author(s)**

Heewon Jeon

## References

Church, K. W. and Hanks, P. (1990). Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22-29.

Church, K. W. and Mercer, R. L. (1993). Introduction to the special issue on computational linguistics using large corpora. *Computational Linguistics*, 19(1):1-24.

---

reloadAllDic	<i>reload all Hannanum analyzer dictionary</i>
--------------	------------------------------------------------

---

## Description

Mainly, user dictionary reloading for Hannanum Analyzer. If you want to update user dictionary on KoNLP\_dic/current/dic\_user.txt, need to execute this function after editing dictionary.

## Usage

```
reloadAllDic()
```

## Examples

```
## Not run:
## This codes can not be run
## if you don't have encoding system which can en/decode
## Hanguk(ex) CP949, EUC-KR, UTF-8).
dicpath <- file.path(system.file(package="Sejong"), "dics", "handic.zip")
conn <- unz(dicpath, file.path("data", "kE", "dic_user2.txt"))
newdic <- read.table(conn, sep="\t", header=FALSE, fileEncoding="UTF-8", stringsAsFactors=FALSE)
mergeUserDic(newdic)
## backup merged new dictionary
backupUsrDic(ask=FALSE)
## restore from backup directory
restoreUsrDic(ask=FALSE)
## reloading new dictionary
reloadAllDic()
## End(Not run)
```

---

restoreUsrDic	<i>use for restoring backup dic_user.txt</i>
---------------	----------------------------------------------

---

## Description

Utility function for restoring dic\_user.txt file to dictionary directory.

## Usage

```
restoreUsrDic(ask = TRUE)
```

**Arguments**

ask                    ask to confirm backup

**Examples**

```
## Not run:
## This codes can not be run
## if you don't have encoding system which can en/decode
## Hanguk(ex) CP949, EUC-KR, UTF-8).
dicpath <- file.path(system.file(package="Sejong"), "dics", "handic.zip")
conn <- unz(dicpath, file.path("data","kE","dic_user2.txt"))
newdic <- read.table(conn, sep="\t", header=FALSE, fileEncoding="UTF-8", stringsAsFactors=FALSE)
mergeUserDic(newdic)
## backup merged new dictionary
backupUsrDic(ask=FALSE)
## restore from backup directory
restoreUsrDic(ask=FALSE)
## reloading new dictionary
reloadAllDic()
## End(Not run)
```

---

SimplePos09

*POS tagging by using 9 KAIST tags*

---

**Description**

Do pos tagging using 9 tags uses Hannanum analyzer. see details in [Hannanum](#). Example will be shown in [github wiki](#).

**Usage**

```
SimplePos09(sentence)
```

**Arguments**

sentence            input

**Value**

results of tagged analysis

**References**

Sangwon Park et al(2010). A Plug-In Component-based Korean Morphological Analyzer

---

`SimplePos22`*POS tagging by using 22 KAIST tags*

---

**Description**

Do POS tagging using 22 tags uses Hannanum analyzer. see details in [Hannanum](#). Example will be shown in [github wiki](#).

**Usage**

```
SimplePos22(sentence)
```

**Arguments**

sentence          input

**Value**

results of tagged analysis

**References**

Sangwon Park et al(2010). A Plug-In Component-based Korean Morphological Analyzer

---

`statDic`*summary of dictionaries*

---

**Description**

show summary, head and tail of current or backup dictionaries

**Usage**

```
statDic(which = "current", n = 6)
```

**Arguments**

which              "current" or "backup" dictionary  
n                    a single integer. Size for the resulting object to view

**Examples**

```
## show current dictionary's summary, head, tail  
statDic("current", 10)
```

---

StoK	<i>Sejong tag to KAIST tag</i>
------	--------------------------------

---

**Description**

Sejong tag to KAIST tag

**Author(s)**

Heewon Jeon <madjakarta@gmail..com>

---

tags	<i>tag names</i>
------	------------------

---

**Description**

tag names

**Author(s)**

Heewon Jeon <madjakarta@gmail..com>

---

useSejongDic	<i>use Sejong noun dictionary</i>
--------------	-----------------------------------

---

**Description**

Retrive Sejong dictionary to use in KoNLP

**Usage**

```
useSejongDic(backup = T)
```

**Arguments**

backup	will backup current dictionary?
--------	---------------------------------

**References**

<http://www.sejong.or.kr/>

---

useSystemDic	<i>use system default dictionary</i>
--------------	--------------------------------------

---

**Description**

Retrive system default dictionary to use in KoNLP

**Usage**

```
useSystemDic(backup = T)
```

**Arguments**

backup	will backup current dictionary?
--------	---------------------------------



# Index

## \*Topic **datasets**

editweights, [6](#)

KtoS, [10](#)

StoK, [15](#)

tags, [15](#)

backupUsrDic, [2](#)

concordance\_file, [3](#)

concordance\_str, [4](#)

convertHangulStringToJamos, [4](#)

convertHangulStringToKeyStrokes, [5](#)

convertTag, [5](#)

editweights, [6](#)

extractNoun, [6](#)

HangulAutomata, [7](#)

is.ascii, [7](#)

is.hangul, [8](#)

is.jaeum, [8](#)

is.jamo, [9](#)

is.moeum, [9](#)

KtoS, [10](#)

mergeUserDic, [10](#)

MorphAnalyzer, [11](#)

mutualinformation, [11](#)

reloadAllDic, [12](#)

restoreUsrDic, [12](#)

SimplePos09, [13](#)

SimplePos22, [14](#)

statDic, [14](#)

StoK, [15](#)

tags, [15](#)

useSejongDic, [15](#)

useSystemDic, [16](#)