

분석프로그래밍 I.11

문자열

국민대학교 경영학부 빅데이터경영통계학 전공

2018. 5. 27(월)

1 인코딩 관련 함수들

- 인코딩 바꾸기
 - 레이블만 바꾸기 `Encoding(x) <- "UTF-8"`
 - 레이블과 내용을 함께 바꾸기 `iconv(x, to="UTF-8")`
- 인코딩 추측하기
 - `stringi::stri_enc_detect(x)`¹
- 인코딩 종류
 - `iconvlist()`
- 인코딩을 다루는 예
 - Listing 1 : R의 문자열은 문자를 특정한 인코딩 방법을 사용하여 저장한다. 인코딩이란 문자를 컴퓨터가 저장하는 방법을 의미한다. 구체적으로 (컴퓨터가 다룰 수 있는) 숫자(0/1 또는 16진수)를 문자 하나에 연결하는 방법을 의미한다. 대표적으로 ASCII 코드에서는 'A'는 16진수로 41, 'B'는 16진수로 42이다. 이렇게 컴퓨터가 저장하고 있는 숫자는 특정한 인코딩 방법에 따라 문자로 해석되는데, 이때 인코딩 방법을 알아야 주어진 숫자를 정확한 문자로 해석할 수 있다. 현재 변수 `x`에 저장되어 있는 문자의 인코딩 방법은 `Encoding(x)`로 확인할 수 있고, `Encoding(x) <- "UTF-8"`을 통해 UTF-8으로 변경할 수도 있다.
 - 내용은 동일하게 보존하면서 인코딩을 바꾸는 경우 : Listing 2

¹화일의 경우에는 `readr::guess_encoding()`을 활용한다. 문자열 벡터의 경우에는 `readr::guess_encoding(charToRaw(x))`를 쓸 수 있다.

```
1 x <- "안녕?"
2
3 Encoding(x)
4 charToRaw(x)
5
6 Encoding(x) <- "UTF-8"
7 print(x)
8 charToRaw(x)
9
10 Encoding(x) <- "unknown"
11 print(x)
12 charToRaw(x)
13
14 Encoding(x) <- "latin1"
15 print(x)
16 charToRaw(x)
17
18 iconvlist()
19
20 Encoding(x) <- "CP949"
21 print(x)
22 charToRaw(x)
23
24 y <- iconv(x, to='UTF-8')
25 print(y)
26 Encoding(y)
27
28 x <- 'A'
29 charToRaw(x)
30 y <- '\x41'
31 print(y)
32 Encoding(y) <- 'UTF-8'
33 print(y)
34
35 x <- '안녕'
36 charToRaw(x)
```

```
37 y <- '\xbe\x8\x3\xe7\xf'
38 print(y)
39 Encoding(y)
```

Listing 1: 문자열 인코딩만 바꾸는 함수, Encoding

```
1 x <- "안녕?"
2
3 Encoding(x)
4 charToRaw(x)
5
6 y <- iconv(x, to='UTF-8')
7 print(y)
8 Encoding(y)
9 charToRaw(y)
10
11 x <- 'ABCD'
12 Encoding(x)
13 charToRaw(x)
14
15 y <- iconv(x, to='UTF-8')
16 print(y)
17 Encoding(y)
18 charToRaw(x)
```

Listing 2: 문자열의 내용은 보존하면서 인코딩만 변환, iconv

2 문자열 관련 함수들

- 표 1 참조.

3 정규표현식Regular Expression

3.1 참고

- `stringr::str_view()`
- <https://regex101.com/>
- <https://regexr.com/>
- <https://regexper.com/>

3.2 대표적인 정규표현식

- 표 2 참조.

3.3 POSIX 브라켓

- 표 3 참조.

표 1: 문자열 관련 함수들

기능	R의 베이스 함수	package:stringr
알파벳 대/소문자	<code>letters[], LETTERS[]</code>	
알파벳 대/소문자 변환	<code>tolower(), toupper()</code>	<code>str_to_title()</code>
문자 수	<code>nchar()</code>	<code>str_length()</code>
문자열 내 공란 제거	<code>trimws()</code>	<code>str_trim()</code>
문자열 연결	<code>paste(), sprintf()</code>	<code>str_c()</code>
문자열 분해	<code>strsplit()</code>	<code>str_split()</code>
문자열의 일부	<code>substring()</code>	<code>str_sub()</code>
특정 패턴의 존재 여부	<code>grepl()</code>	<code>str_detect(), str_count()</code>
특정 패턴의 위치	<code>gregexpr(), regexec()</code>	<code>str_locate()</code>
특정 패턴의 교체	<code>gsub()</code>	<code>str_replace()</code>
특정 패턴의 추출/배제	<code>regmatches(t, gregexpr(p, t))</code>	<code>str_extract(t, p), str_match</code>

3.4 정규표현식을 활용할 수 있는 함수들

- `dir(list.files)`, `apropos`, `find` 등

3.5 Approximate matching(edit-based distance)

- `stringdist::stringdist(a, b, method =)`
- `method =`
 - `'osa'`
 - `'lv'` : delete, insert, substitution
 - `'dl'` : delete, insert, substitution, transposition
 - `'hamming'`
 - `'lcs'` : delete, insert
 - `'gram'`
 - `'cosine'`
 - `'jaccard'`
 - `'jw'`
 - `'soundex'`

3.6 정규표현식을 활용한 몇 가지 예제

```

1 library(stringr)
2 library(readr)
3
4 # 구텐베르크 프로젝트에서 로미오와 줄리엣 전체를 하나의 문자열로 다운 받는다.
5 txt <- read_file('http://www.gutenberg.org/cache/epub/1112/pg1112.txt')
6
7 # 각 장(Scene)의 장소만 추출한다.
8 match <- str_match_all(txt, "Scene[^.]*\\.\\.r\\n([^\\r]*)\\.\\.r")
9 scenes <- match[[1]][,2]
```

```
10 head(scenes)
11
12 # 각 장(Scene)의 장소만 추출한다. 다른 방법.
13 extract <- str_extract_all(txt, "Scene[.]*[.]\r\n(?:[^\r]*)\r")
14 ex <- gsub("Scene [^.]*.\r\n", "", extract[[1]])
15 scenes <- gsub("\r", "", ex)
16 head(scenes)
17
18 # 로미오와 줄리엣을 한 줄을 하나의 문자열로 다운 받는다.
19 txt <- readLines('http://www.gutenberg.org/cache/epub/1112/pg1112.txt',
20                 encoding='UTF-8')
21
22 l <- grep("^Scene", txt)
23 txt[l+1]
```

Listing 3: 정규표현식을 활용한 예제

표 2: 정규표현식

정규표현식	의미
<code>n2ow</code>	<code>n2ow</code>
<code>[abc]</code>	a, b, 또는 c
<code>[a-z]</code>	문자 a에서 z까지
<code>[0-9]</code>	0에서 9까지
<code>[^abc]</code>	a, b, c를 제외한 모든 문자
<code>[\\^], [\\-]</code>	<code>^</code> , <code>-</code> 를 <code>[]</code> 안에서 나타내기
<code>\\d</code>	숫자
<code>\\D</code>	숫자를 제외한 모든 문자
<code>\\w</code>	알파벳, 한글, 숫자 (alphanumeric)
<code>\\W</code>	알파벳, 한글, 숫자를 제외한 모든 문자 (non-alphanumeric)
<code>\\s</code>	공란 (' ', '\\t', '\\r', '\\n')
<code>\\S</code>	공란을 제외한 모든 문자
<code>\\b</code>	단어 경계 (word boundary): w와 W의 연결
<code>\\B</code>	단어 경계가 아닌 경우
<code>.</code>	모든 문자
<code>\\. </code>	.(마침표)
<code>\\[, \\]</code>	[,]
<code>\\</code>	\\
<code>a{m}</code>	a를 m번 반복
<code>[cd]{m,n}</code>	c 또는 d를 m에서 n번의 반복
<code>*</code>	0번 이상의 반복
<code>+</code>	1번 이상의 반복
<code>?</code>	0번 또는 1번
<code>{m,n}?, *?, +?, ??</code>	더 적은 횟수를 선호
<code>\\{, \\}, *, \\+, \\?</code>	{, }, *, +, ?
<code>^ \$</code>	시작과 끝
<code>()</code>	그룹
<code>(ab(cd))</code>	그룹 내의 그룹
<code>(now never)</code>	now 또는 never
<code>\\^ , \\\$, \\(, \\)</code>	<code>^</code> , <code>\$</code> , <code>(</code> , <code>)</code>
<code>\\1, \\2</code>	첫 번째, 두 번째 그룹

표 3: POSIX 브라켓

POSIX 브라켓	의미	다른표현
[:alnum:]	알파벳과 숫자	[a-zA-Z0-9]
[:alpha:]	알파벳	[a-zA-Z]
[:ascii:]	아스키 문자	[\x00-\x7F]
[:blank:]	Space와 Tab	[\t]
[:cntrl:]	제어문자	[\x00-\x1F\x7F]
[:digit:]	숫자	[0-9]
[:graph:]	볼 수 있는 문자(공란, 제어문자 제외 모든 문자)	[\x21-\x7E]
[:lower:]	소문자	[a-z]
[:print:]	볼 수 있는 문자와 공란(제어문자 제외)	[\x20-\x7E]
[:punct:]	문장부호와 기호	[!\"#\$%&'()*+,-./:;<=>?\\[\`{ }]
[:space:]	모든 공란 문자(줄바꿈 포함)	[\t\r\n\v\f]
[:upper:]	대문자	[A-Z]
[:word:]	단어 문자(문자, 숫자, 밑줄)	[A-Za-z0-9_]
[:xdigit:]	16진수 숫자	[A-Fa-f0-9]