

string manipulation in R

There is a substantial amount of **string manipulation** functions in R that are at the disposal of a programmer. The following discusses only a few of these functions and will additionally omit many optional arguments.

The **grep()** function searches for a specified substring pattern in a vector of **x** strings. If **x** has **n** elements—contains **n** strings—the **grep()** function will return a vector of length up to **n**. Each vector element is an index of **x** where the substring pattern of **x[i]** is found.

```
1 > grep(pattern, x, ...) #function syntax
2 > grep("Well", c("Well Bore", "Well Head", "BOP", "Valve"))
3 [1] 1 2
4
5 > grep("well", c("Well Bore", "Well Head", "BOP", "Valve"))
6 integer(0)
```

Note the **pattern** argument in the above example is case-sensitive.

The **nchar()** function returns the length of string **x** as follows:

```
1 nchar(x, type = "chars", allowNA = FALSE, keepNA = NA) #function syntax
2 > nchar("Well Head")
3 [1] 9
```

The **nchar()** function exhibited difficulties dealing with **NA** values and non-character classes; note that the above function accommodates **NA** values and passing a noncharacter class results in a return error.

The **paste()** function concatenates multiple strings, returning a single string:

```
1 > paste("Well", "Head")
2 [1] "Well Head"
3
4 > paste("Well", "Head", sep = "")
5 [1] "WellHead"
6
7 > paste("Well", "Head", sep = ".")
8 [1] "Well.Head"
9
10 > paste("Production", "Well", "Head")
11 [1] "Production Well Head"
```

The **sep =** argument in the **paste()** can be used to specify a character between the strings.

The **sprintf()** function assembles a string from parts in a formatted manner.

```
1 > i <- 8
2 > s <- sprintf("the square of %d is %d", i, i^2)
3 > s
4 [1] "the square of 8 is 64"
```

Concretely, the function name is intended to indicate the **string print** to a **string**, opposed to the screen. The function states to print the string **"the square of"** follows by the decimal value (base 10) of **i**.

The **substr()** function returns the substring in a given character position range **stop:stop** as passed through.

```
1 > substring("Well Head", 6, 9)
2 [1] "Head"
```

The **strsplit()** function splits a string **x** into an R list of substrings based on *another string split* in **x**.

```
1 > strsplit("11-01-2016", split = "-")
2 [[1]]
3 [1] "11" "01" "2016"
```

The **regexpr()** finds the character position of the first instance of the specified **pattern** within the **text**.

```
1 > regexpr("Head", "Well Head")
2 [1] 6
3 attr(,"match.length")
4 [1] 4
5 attr(,"useBytes")
6 [1] TRUE
```

The above illustration reports that “**Head**” appeared in “**Well Head**”, beginning at character position **6**.

The **gregexpr()** function is synonymous to the **regexpr()** with the addition of finding **all** string instances.

```
1 > gregexpr("ing","Cementing and Casing")
2 [[1]]
3 [1] 7 18
4 attr(,"match.length")
5 [1] 3 3
6 attr(,"useBytes")
7 [1] TRUE
```

The above illustration finds “**ing**” twice in “**Cementing and Casing**”, starting at character position 7 and 18.

regular expressions [T] in R

programming languages occasionally reference **regular expression** in the context of **string manipulation**. The latter becomes relevant when using any of the below string functions in R:

The following functions should be carefully used with regular expressions in R:

grep()	grepl()	regexpr()	gregexpr()	sub()	gsub()	strsplit()
--------	---------	-----------	------------	-------	--------	------------

A **regular expression** is shorthand to specify broad classes of strings. For example, the expression “[**We**]” refers to any string the contains **either** of the letters **W** or **e**; illustrated as follows:

```
1 > grep("[We]", c("Production", "Well", "Head"))
2 [1] 2 3
```

The above illustration indicates the vector elements **2** and **3** contain either a “**W**” or a “**e**” (case sensitive).

```
1 > grep("W.l", c("Production", "Well", "Head"))
2 [1] 2
3
4 > grep("w.l", c("Production", "Well", "Head"))           #case sensitive!
5 integer(0)
```

The above illustration indicates the vector element **2** contains a “**W**” followed by any character before an “**l**”.

```
1 > grep("P..d", c("Production", "Well", "Head"))
2 [1] 1
```

The above illustration indicates the vector element **1** contains a “**P**” followed by any 2-characters before “**d**”.

In the case where an actual period “.” is to be searched, the following displays the importance of syntax:

```
1 > grep(".", c("Prod.", "Well", "Hd."))
2 [1] 1 2 3
3
4 > grep("\\.", c("Prod.", "Well", "Hd."))
5 [1] 1 3
```

Note that specifying a period “.” Chooses elements **1**, **2**, and **3**; periods are **metacharacters**.

Consequently, the **metacharacter** nature of the period must be **escaped** prior to calling the function. The latter is achieved by using a backslash “\”; note, **two** are required due to the first backslash needing escape.

string utilities in the **edtdbg** debugging tool ≡ in R

The code within the **edtdbg** debugging tool (discussed later) utilizes string utilities. The **dgbssendeditcmd()** functions is an illustration of the latter:

```
1 > dgbssendeditcmd                                #send command to editor
2 function(cmd) {
3     syscmd <- paste("vim --remote-send ", cmd, " iiservername ", vimserver, sep = "")
4 }
```

The main concept of the above illustration is that the **edtdbg** sends remote commands to the **Vim** text editor. For example, the following instructs **Vim** on server **168** to move the cursor to line **12** (typed into a terminal *shell* window):

```
1 vim --remote-send 12G --servername 168
```

Understanding the above illustration, it is shown how the **dgbssendeditcmd()** function incorporates R programming to deliver the same results:

```
1 > paste("vim --remote-send ",cmd," --servername ",vimserver,sep="")
2 vim --remote-send 12G --servername 168
```

An additional core element in the operation of **edtdbg** is an arrangement to record the file **dbgsink** containing most R debugger data through a call to the **sink()** function. The information captured includes the line numbers and cursor position in the source file as it is stepped through with R's debugger:

```
1 debug at cities.r#16: {
```

The following illustration continues to discuss other core functions of the **editdbg** tool in R that extensively demonstrates the usefulness of many R **string manipulation** functions. Details on the below code snippets can be found on pages 258-259 of ***The Art of R Programming*** by Norman Matloff:

```
1 > linenumstart <- regexpr("#",debugline) + 1
2 > buffname <- substr(debugline,10,linenumstart-2)
3 > colon <- regexpr(":",debugline)
4 > linenum <- substr(debugline,linenumstart,colon-1)
```

```
1 > substr(debugline,10,linenumstart-2)
```

```
1 > kbodin <- readline(prompt="enter number(s) of fns you wish to toggle dbg: ")
2 enter number(s) of fns you wish to toggle dbg:
```

```
1 1 4 5
```

```
1 > tognums <- as.integer(strsplit(kbodin,split=" ")[[1]])
```