

Lecture 14: Regular expressions in R

STAT598z: Intro. to computing for statistics

Vinayak Rao

Department of Statistics, Purdue University

```
In [ ]: options(repr.plot.width=5, repr.plot.height=3)
```

We have seen the print function:

```
In [ ]: x <- 1
print(x)
y <- list('Hello', TRUE, c(1,2,3))
print(y)
```

print is a *generic* function:

- looks at class of input and calls appropriate function

```
In [ ]: my_df <- data.frame(x = c(1,2), y = c(3,4))
print(my_df)
```

```
In [ ]: print.default(my_df)
```

```
In [ ]: print.data.frame(my_df)
```

```
In [ ]: class(df) <- NULL
print(my_df)
```

print and cat

print can only print its first term

```
In [ ]: print('Right now it is', date())
```

For this we need the cat (concatenate) function

```
In [ ]: cat('Right now it is', date(), "in West Lafayette")
```

```
cat(..., file = '' , sep = ' ' , fill = FALSE,  
labels = NULL, append = FALSE)
```

...: Inputs that R concatenates to print

sep: What to append after each input (default is space)

file: Destination file (default is stdout)

Use `paste()` to store the concatenated output (a string)

```
In [ ]: cat(1:5)
```

```
In [ ]: cat(1:5, sep= ' , ' )
```

```
In [ ]: cat(1:5, sep= '\n' )
```

```
In [ ]: cat([' ,1:5, ']' , sep=( ' , ' ))
```

```
In [ ]: cat([' ,1:5, ']' , sep=c(' ' , rep(' ' ,4), ' ' ))
```

```
In [ ]: cat('Hello', 'World', 'New para', sep='\n', file='new_file.txt')
```

Section 8.1.22 in *The R Inferno*, Patrick Burns:

- `print` outputs all characters in the string
- `cat` outputs what the string represents

Compare:

```
In [ ]: print('Hello\nBye')
```

```
In [ ]: cat('Hello\nBye')
```

- `'\'` escapes the following character (indicating it is special)

What if we want to output `'\n'` using `cat` ?

Escape `\` with another `\`

```
In [ ]: cat('Hello\\n')
```

Regular expression: representation of a collection of strings

Useful for searching and replacing patterns in strings

Composed of a grammar to build complicated patterns of strings

R has functions, which coupled with regular expressions allow powerful string manipulation

E.g. `grep`, `grepl`, `regexpr`, `gregexpr`, `sub`, `gsub`

Matching simple patterns

```
In [ ]: cities <- c('lafayette', 'indianapolis', 'cincinnati')
        grep('in', cities)
```

```
In [ ]: grepl('in', cities)
```

Usage:

```
grep(pattern, x, ignore.case = FALSE, perl = FALSE, value = FALSE)
```

```
In [ ]: grep('in', cities, value=TRUE) #Return values instead of indices
```

Where in each element did the match occur?

```
In [ ]: regexpr('in', cities)
```

What if more than one match occurred?

```
In [ ]: gregexpr('in', cities)
```

What if we want to match

- any letter followed by 'n'?
- any vowel followed by 'n'?
- two letters followed by 'n'?
- any number of letters followed by 'n'?

Regular expressions!

- allow us to match much more complicated patterns
- build patterns from a simple vocabulary and grammar

R supports two flavors of regular expressions, we will always use `perl` (set option `perl = TRUE`)

'.' (period) represents any character except empty string ""

```
In [ ]: vec<-c('ct','at', 'cat', 'caat', 'cart', 'dog', 'rat', 'carert', 'bet')
```

```
In [ ]: grep('.at', vec, perl = TRUE)
```

```
In [ ]: grep('..t', vec, perl = TRUE)
```

+ represents one or more occurrences

```
In [ ]: vec<-c('ct','at', 'cat', 'caat', 'cart', 'dog', 'rat', 'carert', 'bet')
```

```
In [ ]: grep('ca+t', vec, perl = TRUE)
```

```
In [ ]: grep('c.+t', vec, perl = TRUE)
```

* represents zero or more occurrences

```
In [ ]: vec<-c('ct','at', 'cat', 'caat', 'cart', 'dog', 'rat', 'carert', 'bet')
```

```
In [ ]: grep('c.*t', vec, perl = TRUE)
```

Group terms with parentheses '(' and ')'

```
In [ ]: vec<-c('ct','at', 'cat', 'caat', 'cart', 'dog', 'rat', 'carert', 'bet')
```

```
In [ ]: grep('c(.r)+t', vec, perl = TRUE)
```

```
In [ ]: grep('c(.r)*t', vec, perl = TRUE)
```

'.', '+', '*' are all metacharacters

Other useful ones include:

- ^ and \$ (start and end of line)

```
In [ ]: grep('e.$', vec, perl = TRUE)
```

| (logical OR)

```
In [ ]: grep('(c.t)|(c.rt)', vec, perl = TRUE)
```

[and] (create special character classes)

[0-7ivx]: any of 0 to 7, i, v, and x

[a-z]: lowercase letters

[a-zA-Z]: any letter

[0-9]: any number

[aeiou]: any vowel

```
In [ ]: grep('[ei]t', vec, perl = TRUE)
```

Inside a character class ^ means "anything except the following characters". E.g.

[^0-9]: anything except a digit

```
In [ ]: grep('[^a]t', vec, perl = TRUE)
```

What if we want to match metacharacters like . or +?

```
In [ ]: vec <- c('ct', 'cat', 'caat', 'caart', 'caaat', 'caaraat',
               'c.t')
         grep('c.t', vec, perl = TRUE) #Is this what we want?
```

Escape them with \

WARNING: a single \ doesn't work. Why?

```
In [ ]: cat('c\t')
```

R thinks \. is a special character like \n.

Use two \s

```
In [ ]: cat('c\\.t')
```

```
In [ ]: grep('c\\.t', vec, perl = TRUE)
```

```
In [ ]: grep('c\\.t', vec, perl = TRUE)
```

To match a \, our pattern must represent \\

```
In [ ]: my_var <- '\\n'
         grep('\\n', my_var)
```

```
In [ ]: my_var <- ('\\')
        grep('\\\\', my_var)
```

Search and replace

The sub function allows search and replacement:

```
In [ ]: vec <- c('ct', 'cat', 'caat', 'caart', 'caaat', 'caaraat', 'c.t')
        sub('a+', 'A', vec, perl = TRUE)
```

sub replaces only first match, gsub replaces all

Use backreferences \1, \2 etc to refer to first, second group etc

```
In [ ]: gsub('(a+)r(a+)', 'b\\1brc\\2c', vec, perl = TRUE)
```

Use \U, \L, \E to make following backreferences upper or lower case or leave unchanged respectively

```
In [ ]: gsub('(a+)r(a+)', '\\U\\1r\\2', vec, perl = TRUE)
```

```
In [ ]: gsub('(a+)r(a+)', '\\U\\1r\\E\\2', vec, perl = TRUE)
```