

String Processing

Working with Text in R

Jesse Lecy



Text is hard to analyze...

because language is fickle.



STOP CLUBBING, BABY SEALS

Once again, punctuation makes all the difference ...

ATTENTION

**Toilet
ONLY**
for
**DISABLED
ELDERLY
PREGNANT
CHILDREN**

THANK YOU for shopping with us!!!

JOIN OUR GRAMMAR PARTY AT [FB.COM/GRAMMAREADE](https://www.facebook.com/GRAMMAREADE)

HUNTERS

**PLEASE USE
CAUTION
WHEN HUNTING
PEDESTRIANS
USING
WALK TRAILS**

strings

STRINGS



functions

Function	Use
grep()	Find a word or phrase (returns the proper string).
grepl()	Find a word or phrase (returns a logical vector).
regexpr()	Find a part of a word or phrase - very flexible.
agrep()	Find an approximate match.
sub()	Replace the first occurrence of a word or phrase.
gsub()	Replace ALL occurrences of a word or phrase.
<hr/>	
paste()	Combine multiple strings into a single string.
strsplit()	Split one string into multiple strings.
substr()	Extract part of a string.

**GREGP: Globally search for a Regular Expression and Print*

REGULAR EXPRESSIONS

Operator	Use
.	matches any single character (wild card for single character)
*	matches 0 or more characters (wild card for any number of characters)
^	start of a string
\$	end of a string
?	match any time a character appears 0 or 1 times
+	match any time a character appears 1 or more times
	OR statement - match either statement given
[]	OR statement - match any of the characters given
[^]	match any characters EXCEPT those given in the list
\	escape character - turns an operator into plain text

```
strings <- c("^ab", "ab", "abc", "abd", "abe", "ab 12", "ab$")
```

```
# match anything that starts with ab followed by any character
```

```
grep("ab.", strings, value = TRUE)
```

```
## [1] "abc" "abd" "abe" "ab 12" "ab$"
```

```
# match abc OR abd OR abe
```

```
grep("ab[c-e]", strings, value = TRUE)
```

```
## [1] "abc" "abd" "abe"
```

```
# match anything that is NOT abc
```

```
grep("ab[~c]", strings, value = TRUE)
```

```
## [1] "abd" "abe" "ab 12" "ab$"
```

Operator	Use
.	matches any single character (wild card for single character)
*	matches 0 or more characters (wild card for any number of characters)
^	start of a string
\$	end of a string
?	match any time a character appears 0 or 1 times
+	match any time a character appears 1 or more times
	OR statement - match either statement given
[]	OR statement - match any of the characters given
[^]	match any characters EXCEPT those given in the list
\	escape character - turns an operator into plain text

QUANTIFIERS

Operator	Use
*	matches at least 0 times.
.	matches only one time
+	matches at least 1 times.
?	matches at most 1 times.
{n}	matches exactly n times.
{n,}	matches at least n times.
{n,m}	matches between n and m times.

```
strings <- c("ht","hot","hoot","hoot")
```

```
# match at least zero times
```

```
grep("h*t", strings, value = TRUE)
```

```
## [1] "ht" "hot" "hoot" "hoot"
```

```
# match ONLY one time
```

```
grep("h.t", strings, value = TRUE)
```

```
## [1] "hot"
```

```
# match exactly n times
```

```
grep("ho{2}t", strings, value = TRUE)
```

```
## [1] "hoot"
```

Operator	Use
*	matches at least 0 times.
.	matches only one time
+	matches at least 1 times.
?	matches at most 1 times.
{n}	matches exactly n times.
{n,}	matches at least n times.
{n,m}	matches between n and m times.

POSITION

Operator	Use
<code>^</code>	matches the start of the STRING.
<code>\$</code>	matches the end of the STRING.
<code>\\b</code>	matches the empty string at either edge of a WORD.
<code>\\B</code>	matches the string provided it is NOT at an edge of a word.

```
strings <- c("abcd", "cdab", "cabd", "c abd")
```

```
# anywhere in the text
```

```
grep("ab", strings, value = TRUE)
```

```
## [1] "abcd" "cdab" "cabd" "c abd"
```

```
# at the beginning of a STRING
```

```
grep("^ab", strings, value = TRUE)
```

```
## [1] "abcd"
```

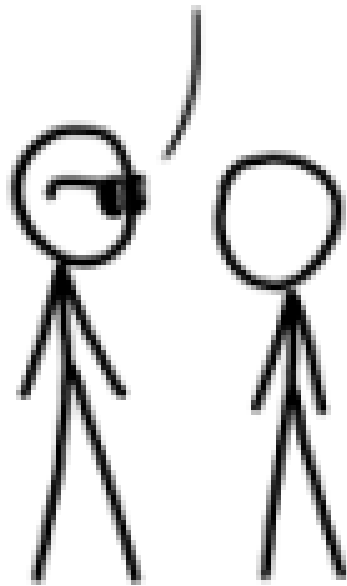
```
# at the end of a STRING
```

```
grep("ab$", strings, value = TRUE)
```

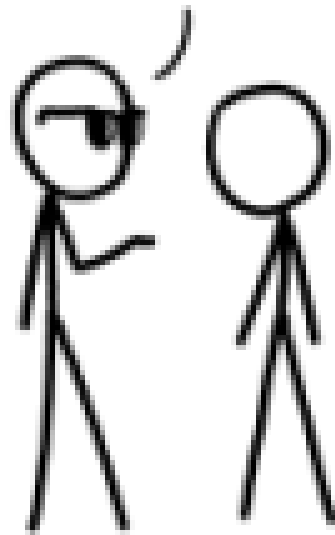
```
## [1] "cdab"
```

Operator	Use
<code>^</code>	matches the start of the STRING.
<code>\$</code>	matches the end of the STRING.
<code>\\b</code>	matches the empty string at either edge of a WORD.
<code>\\B</code>	matches the string provided it is NOT at an edge of a word.

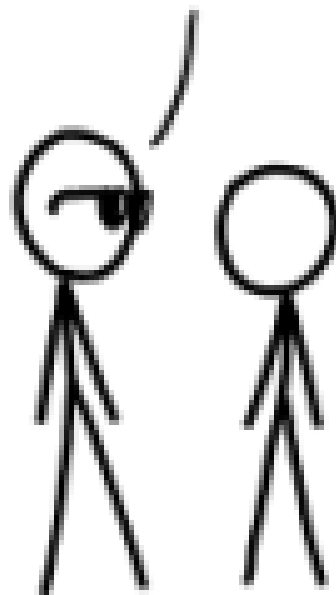
I GOT 99
PROBLEMS,



SO I USED
REGULAR
EXPRESSIONS.



NOW I HAVE
100 PROBLEMS.



Regex in the wild: advanced examples for inspiration

<http://code.tutsplus.com/tutorials/8-regular-expressions-you-should-know--net-6149>

Matching a password:

4. ...and finally the end of the line.

2. ...then six to eighteen...

```
/^[a-z0-9_-]{6,18}$/
```

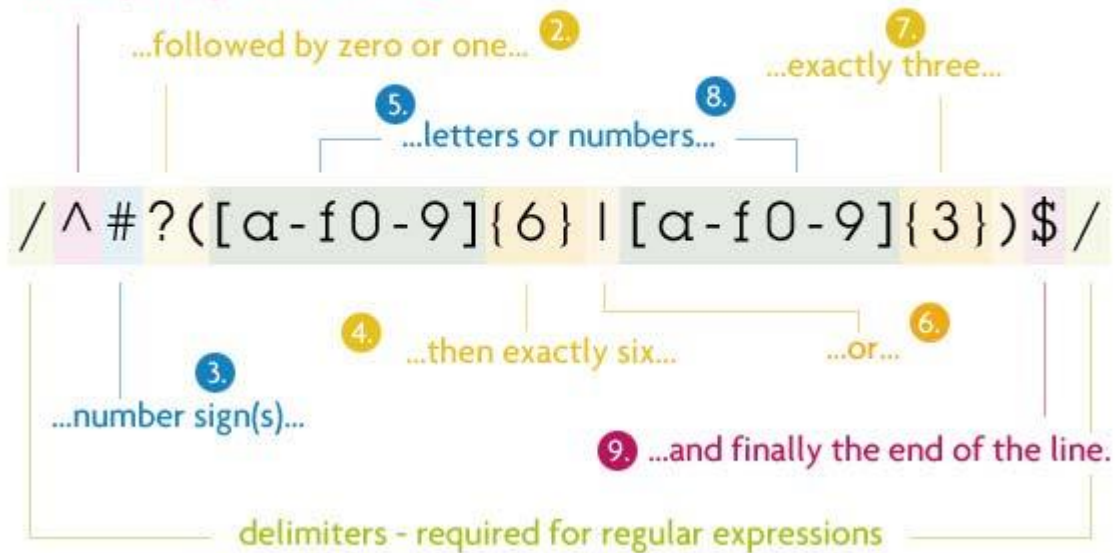
1. The beginning of the line...

3. ...letters, numbers, underscores, or hyphens...

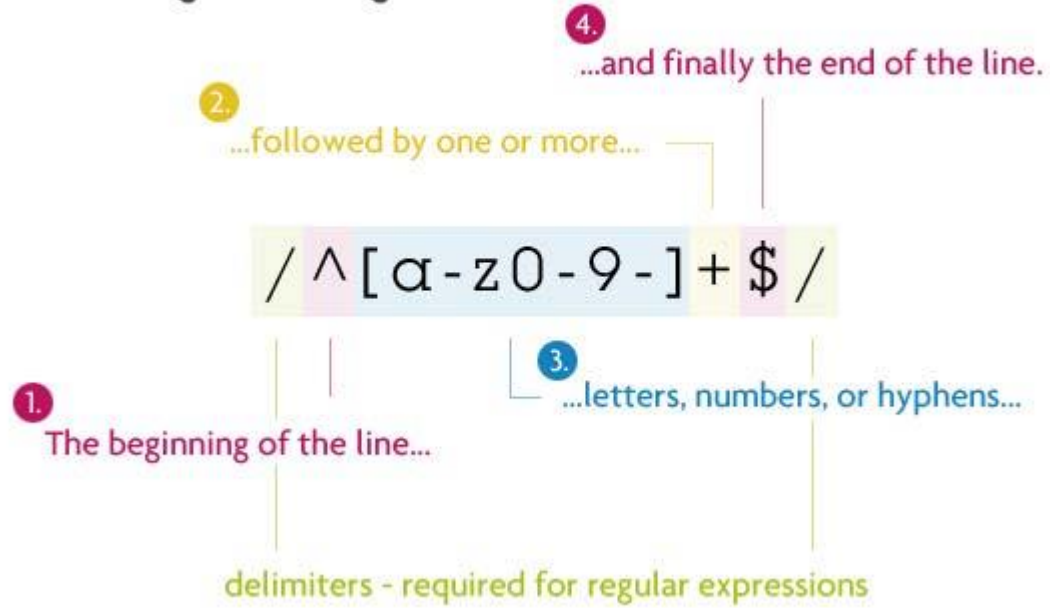
delimiters - required for regular expressions

Matching a hex value:

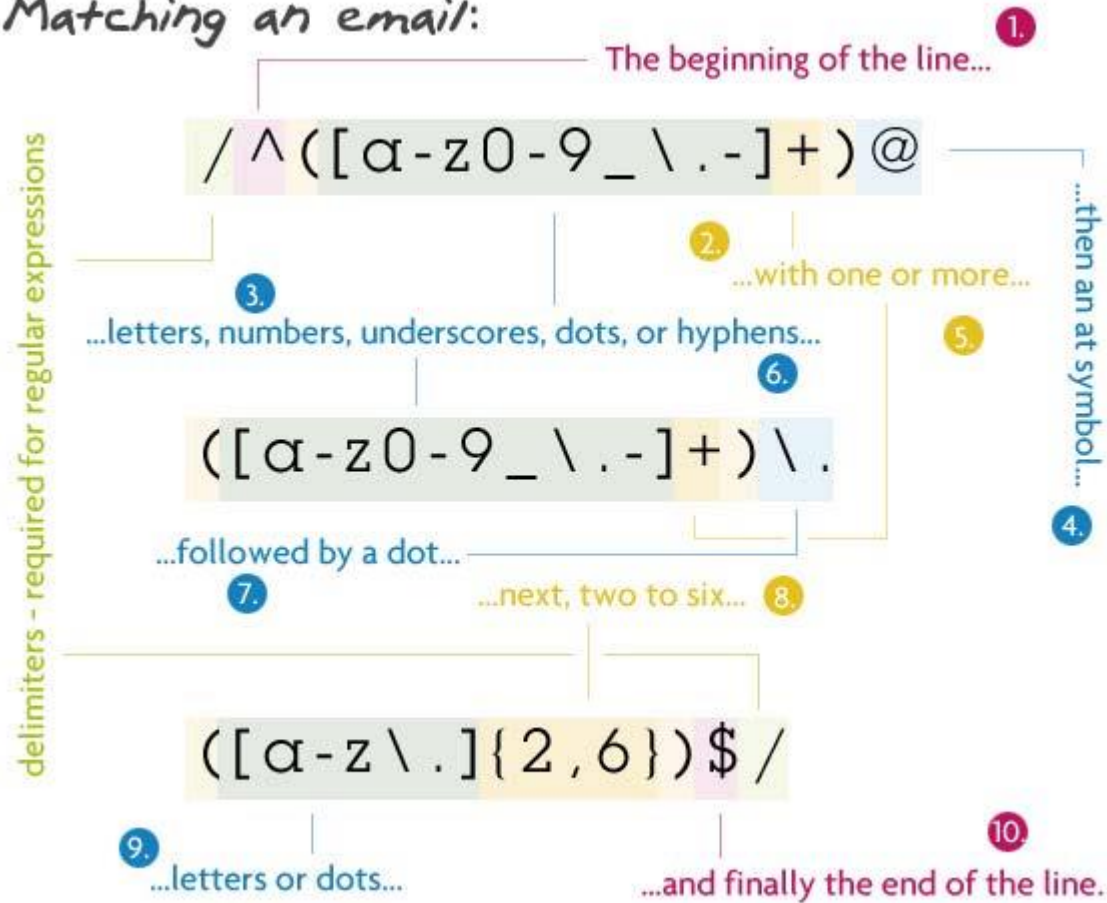
The beginning of the line... 1.



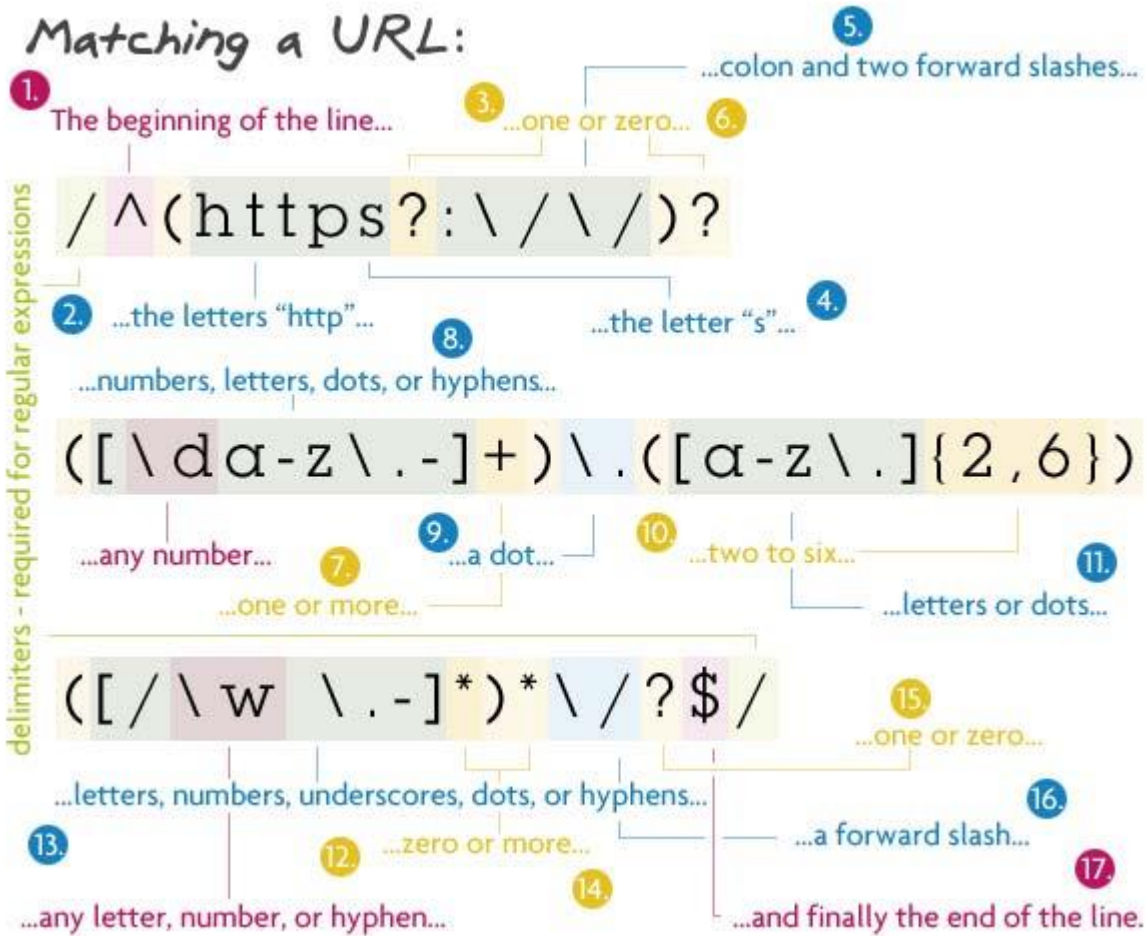
Matching a "slug":



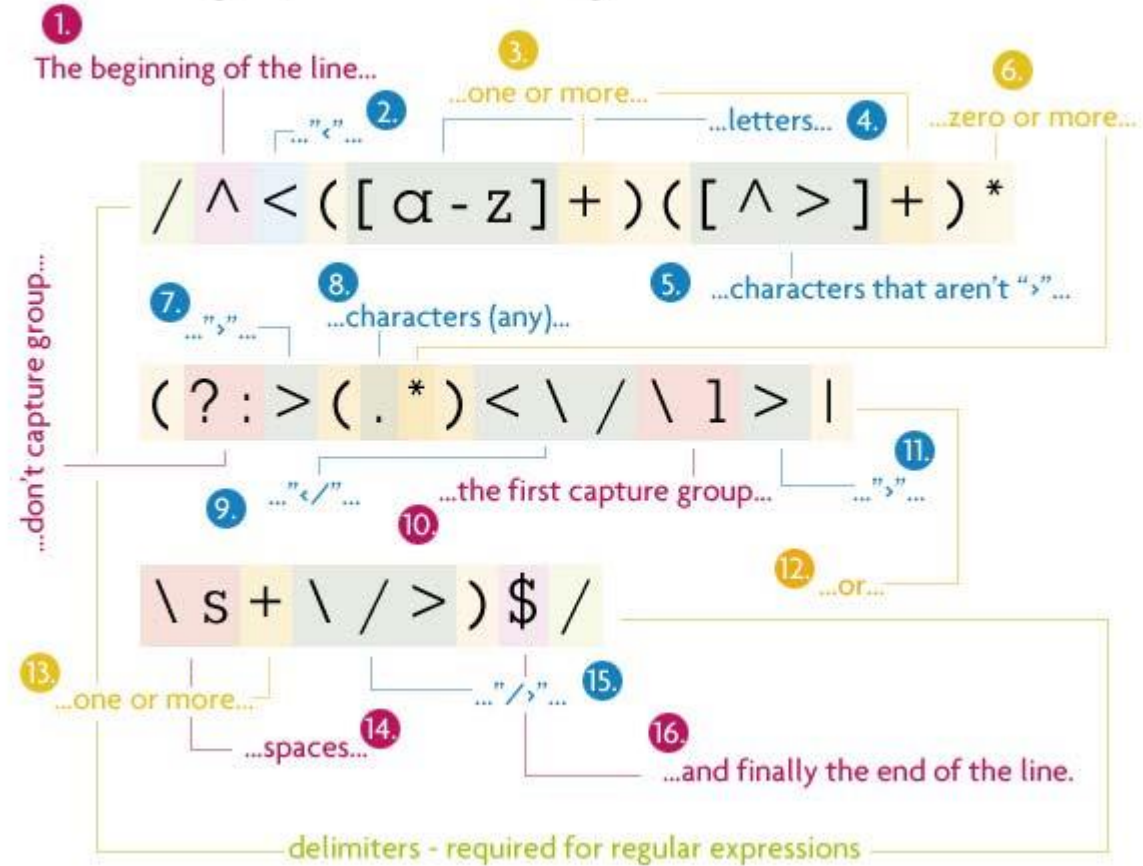
Matching an email:



Matching a URL:



Matching an HTML tag:



fuzzy matching

`agrep()`

ARE TWO WORDS THE SAME?

GEORGE BUSH.....GEORGE W. BUSH

PLANE.....PLAIN

BUREAUCRACY.....BUREACRACY

INTENTION VS. EXECUTION



Minimum Edit Distance

- Two strings and their **alignment**:

I	N	T	E	*	N	T	I	O	N
*	E	X	E	C	U	T	I	O	N



Minimum Edit Distance

I N T E * N T I O N
| | | | | | | | |
* E X E C U T I O N
d s s i s

- If each operation has cost of 1
 - Distance between these is 5
- If substitutions cost 2 (Levenshtein)
 - Distance between them is 8

CALCULATE EDIT DISTANCE

```
> adist( "intention", "execution" )  
      [,1]  
[1,]    5
```

```
> adist( "intention", "execution",  
        costs=c(insertions=1, deletions=1, substitutions=2) )  
      [,1]  
[1,]    8
```

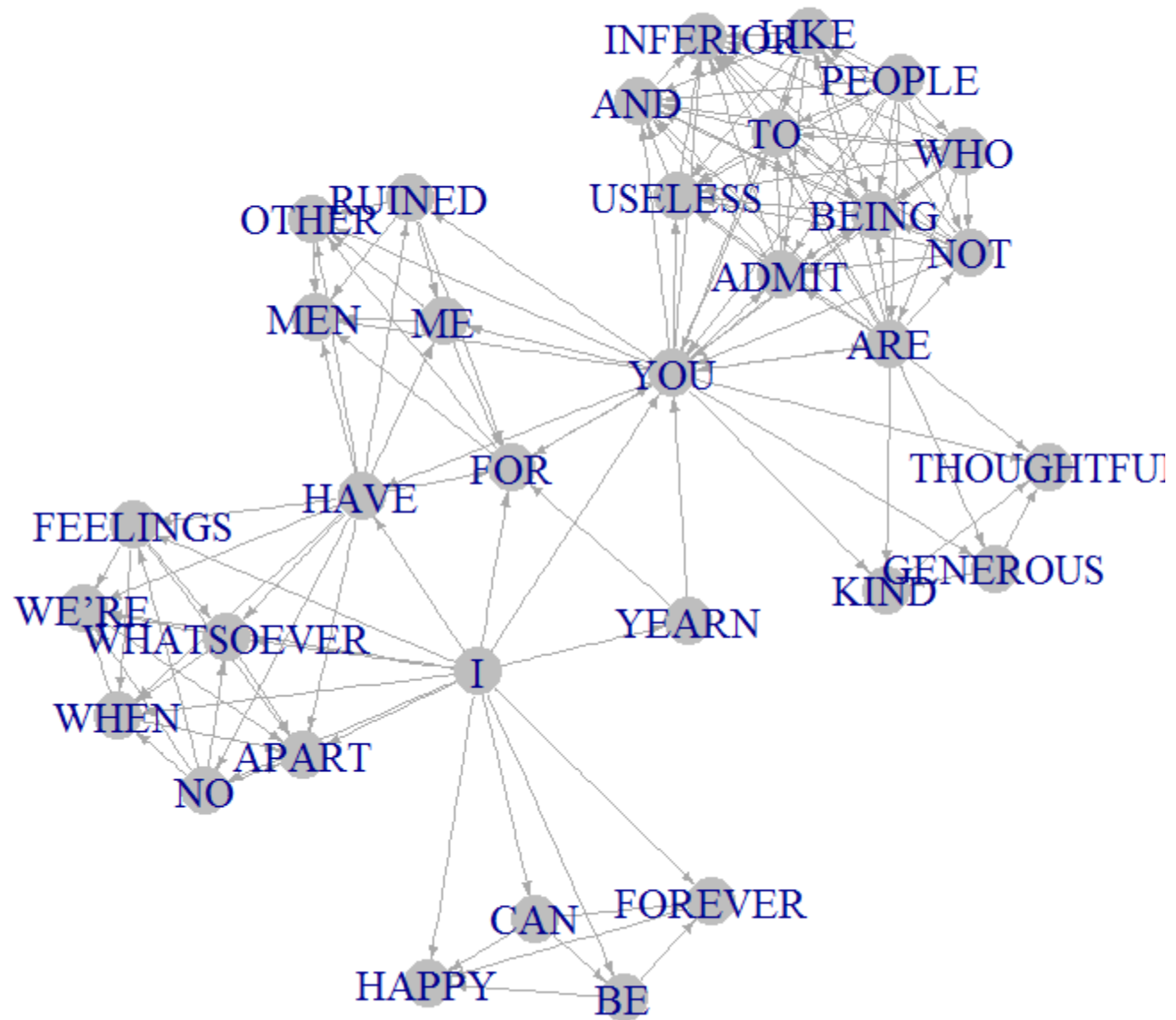
RETURN WORD WITH EDIT DISTANCE < 10%

```
> agrep("lazy", c("daisy", "lasy", "fazed"), value = TRUE )  
[1] "lasy"
```

"grep" stands for "**g**lobally search a **r**egular **e**xpression and **p**rint"

"a" stands for approximate matching

practice



Dear John:

I want a man who knows what love is all about. You are generous, kind, thoughtful. People who are not like you admit to being useless and inferior. You have ruined me for other men. I yearn for you. I have no feelings whatsoever when we're apart. I can be forever happy.

Will you let me be yours?

Gloria

Dear John:

I want a man who knows what love is. All about you are generous, kind, thoughtful people, who are not like you. Admit to being useless and inferior. You have ruined me. For other men, I yearn. For you, I have no feelings whatsoever. When we're apart, I can be forever happy. Will you let me be?

Yours,

Gloria

INSTRUCTIONS

1. Import text
2. Standardize letter case
3. Remove commas, quote marks, special characters
4. Delete empty lines
5. Split the text into sentences
6. Build a network based upon all words in each sentence

FUNCTIONS

Import Text File:

```
readLines( "filename.txt" )
```

Uppercase / Lowercase

```
toupper( x= ), tolower( x= ) # where x is a character vector
```

Find and Replace All:

```
gsub( pattern=, replacement=, x= ) # pattern = what to replace  
# replacement = new text  
# x = character vector
```

Split Text by Delimiter:

```
strsplit( x=, split= ) # x = character vector  
# split = delimiter where splits occur
```

LOAD TEXT FILE

```
setwd( "." )
```

```
x <- readLines("Dear John 1.txt", warn=FALSE)
```

STANDARDIZE CASE

```
x <- toupper( x )
```

REMOVE SPECIAL CHARACTERS

```
x <- gsub( ",", "", x )  
x <- gsub( "we're", "we are", x )  
x <- gsub( "\\:", "", x )  
x <- gsub( "\\?", "", x )
```

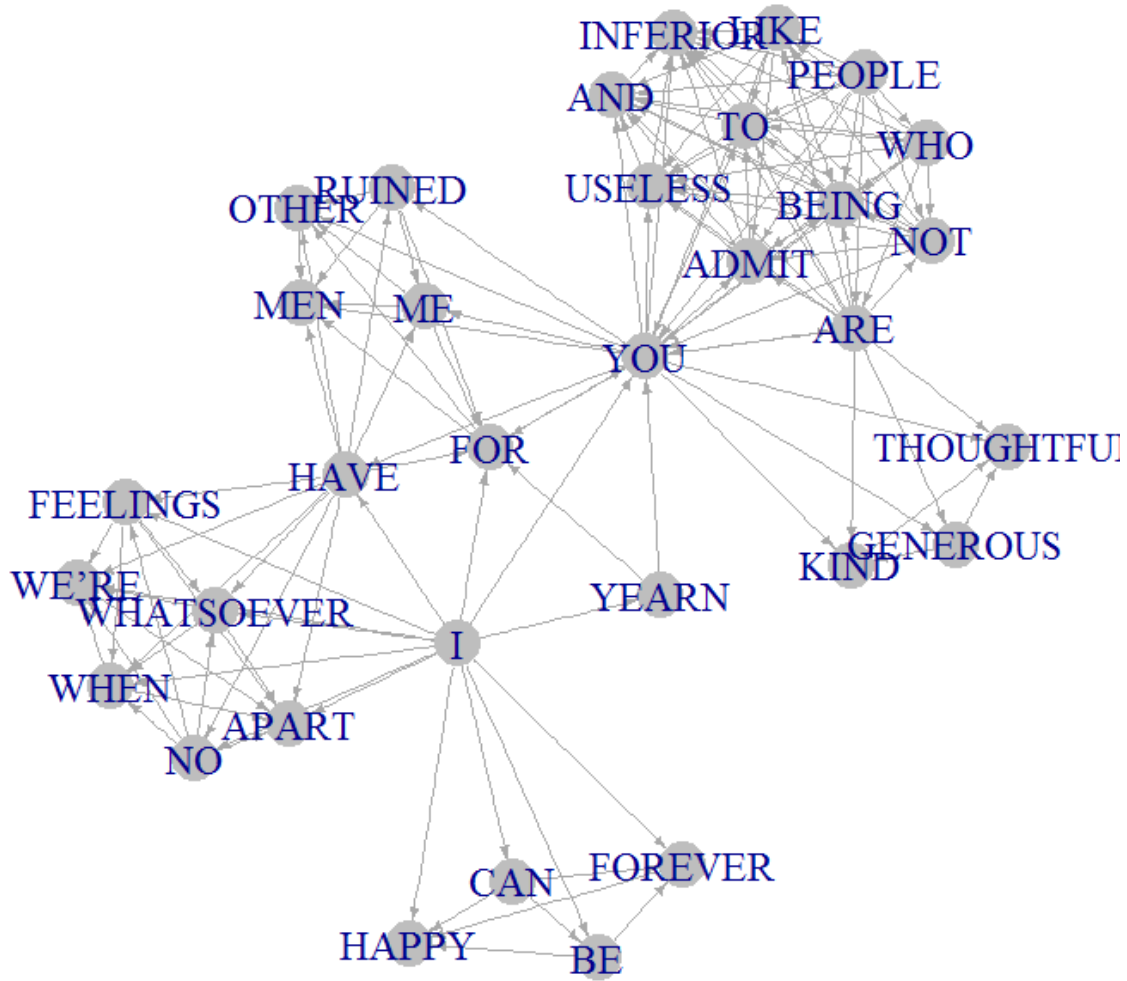
REMOVE BLANK LINES

```
grep( "^$", x )  
x <- x[ - grep( "^$", x ) ]
```

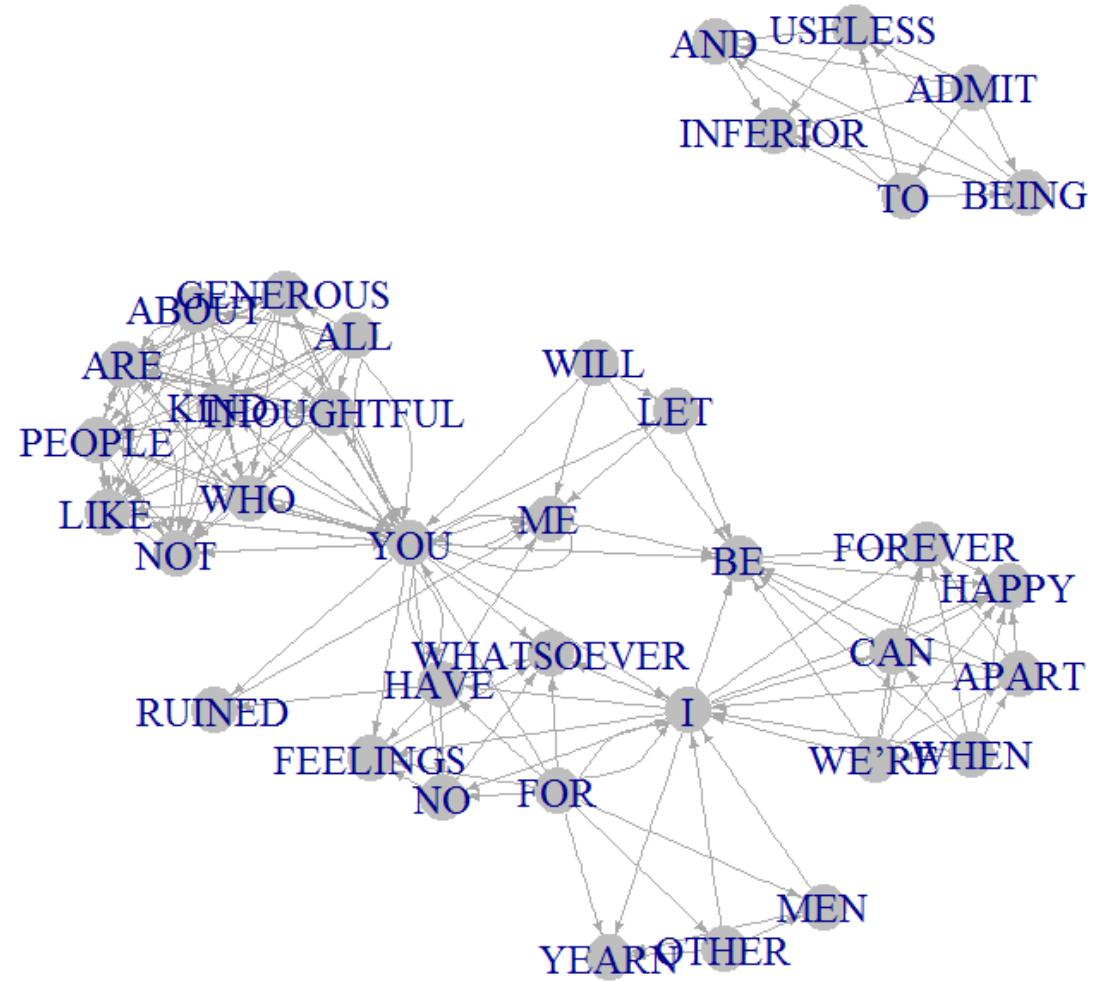

SPLIT INTO SENTENCES

```
strsplit( x, split="\\.")
```

Letter 1



Letter 2



Word Ties Shared in Both Letters

