

Regular Expressions: Disjunctions

- The string of characters inside the braces specifies a disjunction of characters to match

Pattern	Matches
<u>[wW]oodchuck</u>	Woodchuck, woodchuck
<u>[1234567890]</u>	Any digit

[wW]

[ww]

- Ranges [A-Z] the brackets can be used with the dash (-) to specify range any one character in a range.

[A B C D E F ... Z]

Pattern	Matches	
<u>[A-Z]</u>	An upper case letter	<u>D</u> renched Blossoms
<u>[a-z]</u>	A lower case letter	<u>m</u> y beans were impatient
<u>[0-9]</u>	A single digit	Chapter <u>1</u> : Down the Rabbit Hole

(2-5)

[2345]

⊗

(1-9]

[23456789]

Regular Expressions: Negation in Disjunction

Negations [^Ss]

- Carat means negation only when first in []

[^]

Pattern	Matches	
[^A-Z] (^A-Z)	<u>Not an upper case letter</u>	O <u>y</u> fn pripetchik
[^Ss]	Neither ' <u>S</u> ' nor ' <u>s</u> '	<u>I</u> have no exquisite reason"
[e^] (^)	Either e or ^	Look h <u>e</u> re
<u>a</u> ^ <u>b</u>	The pattern a carat b	Look up (<u>a^b</u>) now

[a b] [ww]

Regular Expressions: More Disjunction

^w
Woodchuck is another name for groundhog!

The pipe | for disjunction

Pattern	Matches
groundhog woodchuck	woodchuck
(yours) (mine) [yours mine]	yours
a b c [a b c]	= [abc]
[gG]roundhog(l) [Ww]oodchuck .	Woodchuck



Regular Expressions: Kleene $*$, Kleene $+$

- Kleene $*$: The set of operators that allows us to say things like “some number of as” are based on the asterisk ($*$) commonly called the Kleene*.

$a^* = a, aa, aaa, \dots$

- Kleene $+$ means “one or more occurrences of the immediately preceding character or regular expression”.

- period ($.$), a wildcard expression that matches any single character (except a carriage return).

$a^+ = a, aa, aaa, \dots$
 $a^* = a, aa, aaa, \dots$

Regular Expressions: ? * + .

color | colour

Pattern	Matches	
<u>colou?r</u>	Optional previous char	<u>color</u> <u>colour</u>
<u>oo*h!</u>	0 or more of previous char	<u>oh!</u> <u>ooh!</u> <u>oooh!</u> <u>ooooh!</u> <i>o o o o h</i> <i>!chuhh!</i>
<u>o+h!</u>	1 or more of previous char	<u>oh!</u> <u>poh!</u> <u>oooh!</u> <u>ooooh!</u>
<u>baa+</u>	<i>b a a s, b a a s</i>	<u>baa</u> <u>baaa</u> <u>baaaa</u> <u>baaaaa</u>
<u>beg.n</u>	.	<u>begin</u> <u>begun</u> <u>begun</u> <u>beg3n</u>



Stephen C Kleene

Kleene *, Kleene +

aheen aeha achen

Regular Expressions: Anchors ^ \$

- Anchors are special characters that anchor regular expressions to particular places in a string.

Line start

Pattern	Matches
<u>^</u> [A-Z]	<u>P</u> alo Alto
<u>^</u> [<u>^</u> A-Za-z]	<u>1</u> "Hello"
<u>\.</u> \$	The end <u>.</u>
<u>\.</u> \$	The end <u>?</u> The end <u>!</u>

^
[^ => negate

1 993 99_

1 99 12

\b (word boundary) \B (non word boundary)

Regular Expressions:

Disjunction Operator:

- The disjunction operator, also called the pipe symbol |.
- The pattern /cat|dog/ matches either the string cat or the string dog.

Precedence:

fly | flies fly|flies

- To make the disjunction operator apply only to a specific pattern, we need to use the parenthesis operators (and).

fly (flies)

Example

Find me all instances of the word “the” in a text.

the

The THE

Misses capitalized examples

[tT]he

Incorrectly returns other or theology

[^a-zA-Z] [tT]he [^a-zA-Z]

Errors

RE

the

① ② 3 ④ ⑤
X X X X X

The process we just went through was based on **fixing**
two kinds of errors:

1. Matching strings that we should not have matched (there,
then, other)

False positives (Type I errors) ↘

Increasing precision → minimizing False positive

the
the

2. Not matching things that we should have matched (The)

False negatives (Type II errors)

Increasing recall → minimizing False negatives

More Operators

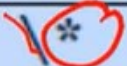


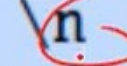

RE	Expansion	Match	First Matches
<u>\d</u>	[0-9]	any digit	Party_of_5
<u>\D</u>	[^0-9]	any non-digit	Blue_moon
<u>\w</u>	<u>[a-zA-Z0-9_]</u>	any <u>alphanumeric/underscore</u>	<u>D</u> aiyu
<u>\W</u>	[^\w]	a non-alphanumeric	<u>!</u> !!!
<u>\s</u>	[\t\r\n\f]	whitespace (space, tab)	
<u>\S</u>	[^\s]	Non-whitespace	<u>i</u> n_Concord

More Operators

RE	Match
*	zero or more occurrences of the previous char or expression
+	<u>one</u> or more occurrences of the previous char or expression
?	exactly zero or one occurrence of the previous char or expression
{n}	<u>n occurrences of the previous char or expression</u>
{n,m}	from <u>n</u> to <u>m</u> occurrences of the previous char or expression
{n,}	at least <u>n</u> occurrences of the previous char or expression
{,m}	up to <u>m</u> occurrences of the previous char or expression

$color^*$ $a\{3\}$ aaa $a\{3,5\}$
 $a \xrightarrow{2} 10$ $(\underline{ab})\{3\}$ $ababab$
 $a\{2,10\}$

More Operators

RE	Match	First Patterns Matched
 *	an asterisk “*”	“K_A*P*L*A*N”
 \.	a period “.”	“Dr_ Livingston, I presume”
 \?	a question mark	“Why don’t they come and lend a hand_?”
 \n	a newline	
 \t	a tab	