

Regular expressions

Regular expressions are a very powerful tool to do string matching and processing

Allows you to do things like:

- □ Tell me if a string starts with a lowercase letter, then is followed by 2 numbers and ends with "ing" or "ion"
- □ Replace all occurrences of one or more spaces with a single
- □ Split up a string based on whitespace or periods or commas
- Give me all parts of the string where a digit is proceeded by a letter and then the '#' sign

WENDER FLISHN A.

NEW DEEL TOROCT
EARGORE FRINGY
SCHOOLS (SEWEN BE DIV.)

THE ROW WORTHOR REVIOUS SOFTWARE FROM ALCOHOLS

SOFTWARE FROM THE LAW ALCOHOLS

TENNESS HE DIV. I KNOW REGULAR EXPRESSIONS. EVERYBODY STAND BACK. http://xkcd.com/208/

Regular expressions: literals

We can put any string in a regular expression

- matches any string that has "test" in it
- /this class/
- matches any string that has "this class" in it
- case sensitive: matches any string that has "Test" in it

4

## Regular expressions: character classes A set of characters to match: put in brackets: [] [abc] matches a single character a or b or c What would the following match? /[Tt]est/ any string with "Test" or "test" in it

5

Regular expressions: character classes

A set of characters to match:

put in brackets: []

[abc] matches a single character a or b or c

Can use - to represent ranges

[a-z] is equivalent to

[A-D] is equivalent to

[0-9] is equivalent to

6

8

# Regular expressions: character classes A set of characters to match: put in brackets: [] [abc] matches a single character a or b or c Can use - to represent ranges [a-z] is equivalent to [abcdefghijklmnopqrstuvwxyz] [A-D] is equivalent to [ABCD] [0-9] is equivalent to [0123456789]

```
Regular expressions: character classes

For example:
/(0-9|(0-9|(0-9)/0.9)/
matches any four digits, e.g. a year

Can also specify a set NOT to match:
^ means all characters EXCEPT those specified

a [^a] all characters except 'a'

a [^0-9] all characters except numbers

a [^A-Z] ???
```

## Regular expressions: character classes For example: /[0-9][0-9][0-9][0-9]/ matches any four digits, e.g. a year Can also specify a set NOT to match: ^ means all characters EXCEPT those specified [^a] all characters except 'a' [^0-9] all characters except numbers [^A-Z] not an upper case letter (be careful, this will match any character that's not uppercase, not just letters

9

11

Meta-characters (not always available)

\[ \w - \word \character \( (a-zA-Z\_0-9) \)

\[ \w - \word \character \( (a-zA-Z\_0-9) \)

\[ \w - \word \character \( (a-zA-z\_0-2) \)

\[ \w - \word \character \

10

## /19\d\d/ would match any 4 digits starting with 19 /\s\s\ matches anything with two adjacent whitespace characters (spaces, tabs, etc) /\s[aeiou]..\s/ any three letter word that starts with a vowel

Regular expressions: repetition

\* matches zero or more of the preceding character
/bord/
matches any string with:

\* bid
\* boad

12

# Regular expressions: repetition ? zero or 1 occurrence of the preceding /fights?/ matches any string with "fight" or "fights" in it {n,m} matches n to m inclusive /ba{3,4}d/ matches any string with = baaad = baaaad

13

15

Regular expressions:
beginning and end

^ marks the beginning of the line
\$ marks the end of the line

/test/ test can occur anywhere

/^test/ must start with test

/test\$/ must end with test

/\*test\$/ \*??

14

Regular expressions:
beginning and end

^ marks the beginning of the line
\$ marks the end of the line

/test/ test can occur anywhere

/^test/ must start with test

/test\$/ must end with test

/\*test\$/ must be exactly test

Regular expressions: repetition revisited

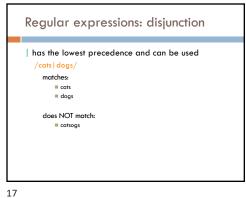
What if we wanted to match:
This is very interesting
This is very very interesting
This is very very interesting
Would /This is very+ interesting/ work?

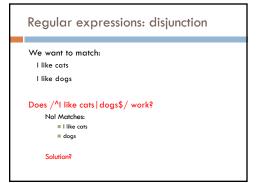
No... + only corresponds to the 'y'

/This is (very )+interesting/

Repetition operators only apply to a single character.
Use parentheses to group a string of characters.

16





18

Regular expressions: disjunction We want to match: I like cats I like dogs /^I like (cats | dogs)\$/ matches: ■ I like cats ■ I like dogs

19

Some examples All strings that start with a capital letter IP addresses 255.255.122.122 Matching a decimal number All strings that end in 'ing' All strings that end in 'ing' or 'ed' All strings that begin and end with the same character

20

### All strings that start with a capital letter /^{A-Z}/ IP addresses /b\d{1,3}\\d{1,3}\\d{1,3}\\d{1,3}\b/ Matching a decimal number /[++|3(0-9)\*,3(0-9)+/ All strings that end in 'ing' /ing\$/ All strings that end in 'ing' or 'ed' /ing | ed\$/

Regular expressions: memory

All strings that begin and end with the same character

Requires us to know what we matched already

()

used for precedence
also records a matched grouping, which can be referenced later

/^(.).\*\1\$/
all strings that begin and end with the same character

21 22

Regular expression: memory

/She likes (\w+) and they like \1/

What would this match?

23

/She likes (\w+) and they like \1/
She likes bananas and they like bananas
She likes movies and they like movies
...

24

### Regular expression: memory

/She likes (\w+) and they like  $\backslash 1/$ 

We can use multiple matches

/She likes (\w+) and (\w+) and they also like  $\1$  and  $\2$ /

### Regular expressions: substitution

Most languages also allow for substitution

s/banana/apple/

substitute first occurrence banana for apple

s/banana/apple/g

substitute all occurrences (globally)

s/^(.\*)\$/\1\1/

duplicate the string, separated by a space

s/\s+/ /g

substitute multiple spaces to a space

25

26

### Regular expressions by language

Java: as part of the String class

String s = "this is a test"

s.matches("test")

s.matches(".\*test.\*")
s.matches("this\\sis .\* test")

s.split(<u>"\\s+</u>")

s.replaceAll(<u>"\\s+"</u>, " ");

Be careful, matches must match the whole string (i.e. an implicit  $^{\Lambda}$  and  $^{\$}$ )

-----

27

### Regular expressions by language

Java: java.util.regex

Full regular expression capabilities

Matcher class: create a matcher and then can use it

String s = "this is a test"

Pattern pattern = Pattern.compile("is\\s+")

Matcher matcher = pattern.matcher(s)

- matcher.matches()
- matcher.find()
- matcher.replaceAll("blah")
- matcher.group()

28

# Regular expressions by language Python: import re s = "this is a test" p = re.compile("test") p.match(s) p = re.compile(".\*test.\*") re.split('\s+', s) re.sub('\s+', ', s)

29

32

grep
| command-line tool for regular expressions (general regular expression print/parser)
| returns all lines that match a regular expression
| grep "@" twitter.posts
| grep "http:" twitter.posts
| can't use metacharacters (\d, \w), use [] instead
| Often want to use "grep -E" (for extended syntax)

31

### sed another command-line tool that uses regular expressions to print and manipulate strings very powerful, though we'll just play with it Most common is substitution: sed "s/ is a / is not a /g" twitter.posts sed "s/ \*//g" twitter.posts sed doesn't have +, but does have \* Can also do things like delete all that match, etc.

Regular expression resources

General regular expressions:

Ch 2.1 of the book

thtp://www.teadar-expressions.info/

good general twaterials
many language specific examples as well

Java

thtp://download.oracle.com/iavase/tutorial/essential/recex/
See also the documentation for java.util.regex

Python

http://doss.ovthon.ora/houto/regex.html

33

