

COMP233-04B Internet Applications

Lecture 19: Regular Expressions

Matching Substrings

Expression

`mall`

matches the string

`smallest`

because it contains a matching string.

To avoid this, use `^` and `$`:

`^mall$`

Regular Expression Characters

<code>a</code>	Character matches itself
<code>.</code>	Match any character
<code>[a-z]</code>	Match any of these characters
<code>[^ABC]</code>	Match any character except these
<code>^</code>	Match to start of string
<code>\$</code>	Match to end of string

Whitespace

- ◆ Users often enter whitespace at the beginning and end of their input.
- ◆ We generally want to ignore this.

`^ *small *$`

Examples

`[mM]ule`

matches:

`mule, Mule`

`[mM].le`

matches:

`male, Mole, mule, MMle, ...`

Another Example

- ◆ Allow the user to type a sentence with a two-digit number (10-99) in it somewhere.

`[1-9][0-9]`

- ◆ How can we find out which number was found in the input?
- ◆ JavaScript returns it as result!

Result of the `match()` Method

```
var text = document.myform.mytext.value;
var regstr = document.myform.regexp.value;
var regexp = new RegExp(regstr);
var result = text.match(regexp);
if (result) {
    // matches!
    document.bgColor = "lightgreen";
    document.myform.matched.value =
        result[0];
}
```

More Characters

<code>\d</code>	A digit (0...9)
<code>\D</code>	Anything but a digit
<code>\w</code>	A word character (letter, digit, ...)
<code>\W</code>	Anything but a word character
<code>\s</code>	Whitespace (space or tab)
<code>\S</code>	Anything but whitespace

Combining Regular Expressions

<code>(...)</code>	Grouping
<code>*</code>	Zero or more times
<code>+</code>	One or more times
<code>?</code>	Zero or one times (optional)
<code>{5}</code>	Exactly 5 times
<code>{3,6}</code>	3-6 times

Applied after subexpressions,
e.g., `[aA]*`.

Escaping Special Characters

Problem:

How to match strings containing
special characters like

`. ^ $ [] \`

Solution:

Escape them with `\`

`[a-z]+\.`**txt**

More Examples

Any word of letters:

`[a-zA-Z]*`

Any nonempty word of letters:

`[a-zA-Z]+`

Any nonempty word of letters that
starts with a capital:

`[A-Z][a-z]*`

Final Note

Regular Expressions can be
combined with OR,

e.g.,

`large|medium|small`
`[a-z]*|[0-9]*`

But ...

there is no AND.