

Reguljära Uttryck

Pontus Haglund

Institutionen för datavetenskap (IDA)

Anpassat från material ursprungligen av: Eric Elfving

3 oktober 2018

- Introduktion
- Syntax
- Användning i Python

- Reguljära uttryck används för att matcha ett mönster (pattern) mot en text.
- Kallas ofta regex eller regexp (från regular expression)
- Skapades på 50-talet och finns tillgängligt i de flesta språk.
- Ett väldigt bra verktyg som kan användas till mycket.

Ett tecken matchar sig själv

Mönster: | a | han | kul
Söksträng: | a anna | han hon hanna | kul kulle ful

Med hakparenteser ([]) skapas en teckenklass

En teckenklass matchar ett av tecknen innanför hakarna

Mönster: a[ab]c

Söksträng: abc aac acabca

Teckenklasser

Man kan även sätta ett intervall inom hakparenteser

Mönster: a[a-c]c

Söksträng: abc aac acc

Med cirkumflex (^) tar du inversen av en teckenklass

Man matchar då alla tecken utom den givna teckenklassen

Mönster: a[^ab]c

Söksträng: adc aqc aaca

Med lodstreck (|) anger du ett logiskt ELLER

Mönster: | a | c | aa | ab

Söksträng: abc | acaab abc

Vanliga parenteser skapar en "grupp"

Grupper kan användas för att ändra operatorers räckvidd och prioritering

Mönster: $a(aa|ab)c$

Söksträng: aabc aaac acaabca

Upprepningar

Följande operatorer skrivs direkt efter det som ska upprepas.

- ? Matchar 0 eller 1
- * Matchar 0 eller flera
- + Matchar minst 1

Mönster: colou?r a[nl]+a S(oe?|ö)der

Söksträng: color colour anna alla Söder Soder

Upprepningar

- Med klammerparenteser ({ }) upprepar man saker ett visst antal gånger.
- $x\{n, m\}$ matchar x upprepat n till m gånger.
- Precis som i Python kan man utelämna m

Mönster: | [0-9]{2,3} | [a-zA-Z]{5} | [a-zA-Z]{4,}

Söksträng: | 17a 5262 | Linkoping | Linköpings Universitet

Speciella teckensekvenser

Teckensekvens	Motsvarighet
\w	[A-Za-z0-9_]
\W	[^A-Za-z0-9_]
\s	Vita tecken
\S	Inte vita tecken
\d	[0-9]
\D	[^0-9]

Mönster: \d+

Söksträng: 123a

Bakåtreferenser

Alla grupper tilldelas ett nummer (indexeras från 0). Med \N refererar man tillbaka till grupp nummer N.

Mönster: (\d)a\0

Söksträng: 0a0 1a6 4a4

Några speciella tecken

Följande tecken har speciell betydelse:

Tecken Matchar

.

Valfritt tecken (utom newline)

^

Början av rad

\$

Slut på rad eller fil

Mönster: | ^ka..e\$ | ^ka..e\$

Söksträng: | kalle | kalle svensson

Escapesekvenser

För att matcha någon av de speciella tecken (metatecken) som nämnts tidigare ({} [] () ^\$. | *+? \) används ett backslash:

Mönster:		\d	\.	\d		\	\a		\{	\w	+	\}
Söksträng:		1	.	2		4	a		{	hej	}	

- Python har, med stöd av sitt stora bibliotek, så klart stöd för regex.
- Importera bara modulen `re`!

<http://docs.python.org/3/library/re.html>

- re har tre vanligt använda funktioner (och så klart fler):
 - `re.match(pat, str)` Söker efter mönstret pat i början av strängen str
 - `re.search(pat, str)` Söker efter första förekomsten av mönstret pat i strängen str
 - `re.findall(pat, str)` Hittar alla förekomster av mönstret pat i strängen str
- De två första ger ett "match object" vid matchning, annars None. `.findall` ger en lista av matchande strängar.

```
import re
val = input('Mata in ett heltal: ')
while re.match('^[0-9]+$', val) is None:
    val = input('Felaktig inmatning, nytt heltal: ')

val = int(val)
```

Det finns ett problem med reguljära uttryck... backslash!

Backslash används ju redan för escapesekvenser i Python. Tänk er att vi ska söka efter texten \begin. Då måste vi skriva såhär:

```
import re  
match = re.search('\\\\\\begin', str)
```

Vi vill söka efter ett \, då behövs ett extra enligt regex-syntax samt ett extra för varje backslash för att inte Python ska tolka det som en escapesekvens!

Detta problem lösas lättast med "råasträngar (raw string)

```
import re
match = re.search(r'\\begin', str)
```

Match objects

Grupper i mönstret kan tas fram i efterhand med funktionen `group`:

```
import re
match = re.search(r'(\d+) (\w+)', str)
if match is not None:
    print(match.group(0)) # Skriver ut hela matchningen
    print(match.group(1)) # Skriver ut första gruppen
    print(match.group(2))
```

Om en grupp var valbar (hade ett frågetecken) kan den ha värdet `None`.

Användning i Python

22/26

findall

```
>>> str="""En sträng som ibland har heltal i sig: 12345
... Den kan vara lite 2313 olika lång. Denna sträng
... är 3 rader lång!"""

>>> lst = re.findall(r'\d+', str)
['12345', '2313', '3']

>>> [int(x) for x in lst]
[12345, 2313, 3]

>>> sum(map(int, lst))
14664
```

Namngivna grupper

I Python kan man namnge sina grupper:

```
>>> import re
>>> text = "Namn: Anna
... Ålder: 32"
>>> for m in re.finditer(r'(?P<key>\w+): (?P<val>\w+)', str):
...     print(m.group('key'), m.group('val'), sep=' => ')
...
Namn => Anna
Ålder => 32
```

http://en.wikipedia.org/wiki/Regular_expression
Wikipedia

<http://www.regular-expressions.info/> En samling
med tutorials, verktyg och guider

<https://docs.python.org/3/howto/regex.html> Mer
om regex i Python

<http://www.regexr.com/> En bra testare av dina uttryck

1. Ett LiU-ID skrivas på formatet FEN där F är de tre första bokstäverna i förnamnet, E de två första i efternamnet och N är tal med två eller tre siffror (beroende på om det är en anställd eller student). Om förnamnet eller efternamnet endast är två tecken långt blir den delen endast två tecken. Skriv ett mönster för att fånga ett LiU-ID.
2. Ett datum i ISO-format skrivas YYYY-MM-DD. Skriv ett mönster för att fånga detta.
3. Utöka ditt regexp för att kontrollera datumets giltighet (full kontroll går inte då vi saknar logik men t.ex. är en månad som inleds med 9 felaktig).
4. Skriv ett Python-program som (med hjälp av reguljära uttryck) hittar och skriver ut alla html-element i en given fil (på hårddisken eller url, välj själv).

www.liu.se