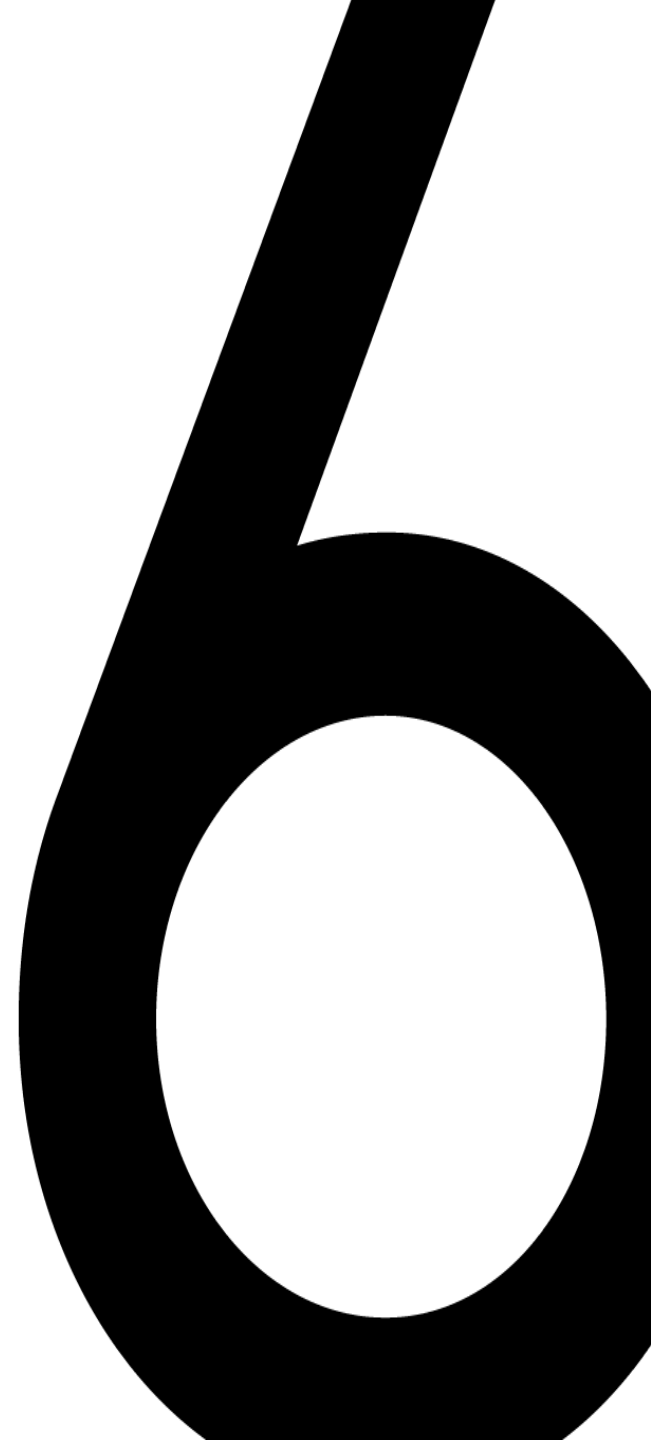

Bash Regex



REGEX

- soorten
- expressies
- =~
- groups
- extglob

Soorten regex

■ POSIX2.0 Regex

- ◆ Extended Regex (ERE)

- `egrep/bash`

- ◆ Oud: Basic Regex(BRE)

- escapen van ? + { | (en)
- standaard grep (zie CS1...)

■ Perl Compatible Regex (PCRE)

- UTF-8 en UNICODE support
- perl/python/java/javascript/c#
- grep -P

Extended en basic regex voorbeeld

■ Extended regex met grep (egrep):

- ◆ `grep -E "groen|rood" bestand.txt`

■ Basic regex met grep

- ◆ `grep "groen\\|rood" bestand.txt`

Enkele regex expressies

■ .

■ [a-z]

■ [0-9]

■ [a|b]

■ a{4}

■ a{1,4}

■ +

■ *

Enkele regex expressies

- . 1 karakter
- [a-z] karakter a tot en met z in kleine letters
- [0-9] getal 0 tot en met 9
- [a|b] karakter a of karakter b
- a{4} vier keer a
- a{1,4} van één tot 4 keer a
- + 1 tot n keer
- * 0 tot n keer

POSIX: Soorten karakters extended regex

- `[[:digit:]]` een cijfer
- `[[:space:]]` spatie, tab, newl, return
- `[[:alnum:]]` letters en cijfers
- `[[:alpha:]]` letters
- `[[:blank:]]` spatie en tab
- `[[:lower:]]` lowercase letters
- `[[:print:]]` afdrukbare karakters

= ~ operator

```
#!/bin/bash
```

```
content="Karel de Grote-Hogeschool, Nationalestraat 5,B-  
2000 Antwerpen"
```

```
regex="B-[0-9]{4}"
```

```
if [[ $content =~ $regex ]]
```

```
then
```

```
    echo "postnummer gevonden" ; exit 0
```

```
else
```

```
    echo "postnummer niet gevonden" >&2 ; exit 1
```

```
fi
```

OPGELET: \$regex is NOOIT met dubbele quotes!

`${BASH_REMATCH[0]}`

- `${BASH_REMATCH[0]}` is een array met het resultaat van de match. Element 0 is de hele match

```
#!/bin/bash
```

```
# Functie: Regex voor een postnummer
```

```
content="Karel de Grote-Hogeschool, Nationalestraat 5,B-  
2000 Antwerpen"
```

```
regex="B-[0-9]{4}"
```

```
[[ $content =~ $regex ]]
```

```
echo "${BASH_REMATCH[0]}"
```

Output: B-2000

regex voorbeeld: email adres

jan . celis @ kdg . be

$[[[:\text{alnum:}]]+ \cdot [[[:\text{alnum:}]]+ @ [[[:\text{alnum:}]]+ \cdot [[[:\text{alpha:}]]+$

Het eerste stuk kan ook zonder punt zijn, dus dat plaatsen we optioneel:

jancelis @ kdg . be

$(([[[:\text{alnum:}]]+ \cdot)\{0,2\} [[[:\text{alnum:}]]+ @ [[[:\text{alnum:}]]+ \cdot [[[:\text{alpha:}]]+$

Het tweede deel kan ook uit subdomeinen bestaan:

jan . celis @ student . kdg . be

$(([[[:\text{alnum:}]]+ \cdot)\{0,2\} [[[:\text{alnum:}]]+ @ (([[[:\text{alnum:}]]+ \cdot)\{1,3\} [[[:\text{alpha:}]]+$

Het laatste domein deel is van 2 tot 3 karakters (zonder cijfers):

jan . celis @ student . kdg . be

$(([[[:\text{alnum:}]]+ \cdot)\{0,2\} [[[:\text{alnum:}]]+ @ (([[[:\text{alnum:}]]+ \cdot)\{1,3\} [[[:\text{alpha:}]]\{2,3\}$

regex groups

- Een group selecteert met haakjes een onderdeel van de regex
- Het resultaat zit in array `n` `${BASH_REMATCH[n]}`

```
content="Karel de Grote-Hogeschool, Nationalestraat 5,B-2000  
Antwerpen"
```

```
regex="([- a-zA-Z]+), ([[:alpha:]]+) ([[:digit:]]+),(.*) (.*)"
```

```
[[ $content =~ $regex ]]
```

```
echo "Naam: ${BASH_REMATCH[1]}"
```

```
echo "Straat: ${BASH_REMATCH[2]} Nr: ${BASH_REMATCH[3]}"
```

```
echo "Postcode: ${BASH_REMATCH[4]} Stad:  
${BASH_REMATCH[5]}"
```

Output:

Naam: Karel de Grote-Hogeschool

Straat: Nationalestraat Nr: 5

Postcode: 2000 Stad: Antwerpen

"Regex" met parameter substitution

```
#!/bin/bash
```

```
content=" Karel de Grote-Hogeschool, Nationalestraat 5,B-2000  
Antwerpen"
```

```
shopt -q -s extglob
```

```
content2=${content#* + ([[digit:]])}  
content=${content##* + ([[digit:]])}
```

```
shopt -q -u extglob
```

```
echo "ungreedy: $content2"
```

```
echo "greedy: $content"
```

Output:

```
ungreedy: ,B-2000 Antwerpen
```

```
greedy: Antwerpen
```

Linux

- Vraag:
 - Ik begrijp de opbouw van de regex niet helemaal die gebruikt wordt in de parameterssubstitutie in volgende voorbeelden:
 - 1) zie cursus p 42:
`content=${content##*+([[:digit:]])}` Dit verwijdert alles tot en met het eerste cijfer. Wat betekent de combinatie van * en + hier?
 - 2) `content=${content##*([[:space:]])}` Dit verwijdert alle leading spaties. Ook hier begrijp ik niet goed wat het sterretje wil zeggen.

Linux

- Verklaring:
 - Extglob = extended globbing
 - Zoals regular expressions, maar met een paar eigenaardigheden
 - quantifiers (+ * ? ...) hebben dezelfde betekenis als in regex, maar moeten voor het patroon (dat tussen ronde haakjes staat) staan ipv erachter

Linux

- Verklaring:
 - <http://wiki.bash-hackers.org/syntax/pattern>
 - Let op: regular expressions belangrijker! (Slides als randinfo.)

Linux

- Verklaring:

Normal pattern language

Sequence	Description
<code>*</code>	Matches any string , including the null string (empty string)
<code>?</code>	Matches any single character
<code>x</code>	Matches the character <code>x</code> which can be any character that has no special meaning
<code>\x</code>	Matches the character <code>x</code> , where the character's special meaning is stripped by the backslash
<code>\\</code>	Matches a backslash
<code>[...]</code>	Defines a pattern bracket expression (see below). Matches any of the enclosed characters at this position.

Linux

- Verklaring:

Bracket expressions



The bracket expression `[...]` mentioned above has some useful applications:

Bracket expression	Description
<code>[XYZ]</code>	The "normal" bracket expression, matching either <code>X</code> , <code>Y</code> or <code>Z</code>
<code>[X-Z]</code>	A range expression: Matching all the characters from <code>X</code> to <code>Y</code> (your current locale , defines how the characters are sorted!)
<code>[:class:]</code>	Matches all the characters defined by a POSIX® character class: <code>alnum</code> , <code>alpha</code> , <code>ascii</code> , <code>blank</code> , <code>cntrl</code> , <code>digit</code> , <code>graph</code> , <code>lower</code> , <code>print</code> , <code>punct</code> , <code>space</code> , <code>upper</code> , <code>word</code> and <code>xdigit</code>
<code>[^...]</code>	A negating expression: It matches all the characters that are not in the bracket expression
<code>[!...]</code>	Equivalent to <code>[^...]</code>
<code>[...] or [-...]</code>	Used to include the characters <code>]</code> and <code>-</code> into the set, they need to be the first characters after the opening bracket
<code>[=C=]</code>	Matches any character that is equivalent to the collation weight of <code>C</code> (current locale!)
<code>[.SYMBOL.]</code>	Matches the collating symbol <code>SYMBOL</code>

Linux

- Verklaring:

Examples

Some simple examples using normal pattern matching:

- Pattern `"Hello world"` matches
 - `Hello world`
- Pattern `[Hh]"ello world"` matches
 - `⇒ Hello world`
 - `⇒ hello world`
- Pattern `Hello*` matches (for example)
 - `⇒ Hello world`
 - `⇒ Helloworld`
 - `⇒ HelloWoRld`
 - `⇒ Hello`
- Pattern `Hello world[:punct:]` matches (for example)
 - `⇒ Hello world!`
 - `⇒ Hello world.`
 - `⇒ Hello world+`
 - `⇒ Hello world?`
- Pattern `[[.backslash.]]Hello[[.vertical-line.]]world[[.exclamation-mark.]]` matches (using collation sybols)
 - `⇒ \Hello|world!`

Linux

- Verklaring:

Extended pattern language

If you set the [shell option](#) `extglob`, Bash understands some powerful patterns. A `<PATTERN-LIST>` is one or more patterns, separated by the pipe-symbol (`PATTERN|PATTERN`).

<code>?(<PATTERN-LIST>)</code>	Matches zero or one occurrence of the given patterns
<code>*(<PATTERN-LIST>)</code>	Matches zero or more occurrences of the given patterns
<code>+(<PATTERN-LIST>)</code>	Matches one or more occurrences of the given patterns
<code>@(<PATTERN-LIST>)</code>	Matches one of the given patterns
<code>!(<PATTERN-LIST>)</code>	Matches anything except one of the given patterns

Linux

- Verklaring:
 - Extglob = extended globbing
 - Zoals regular expressions, maar met een paar eigenaardigheden
 - quantifiers (+ * ? ...) hebben dezelfde betekenis als in regex, maar moeten voor het patroon (dat tussen ronde haakjes staat) staan ipv erachter

Linux

- Vraag:
 - 1) `content=${content##* + ([:digit:])}`
 - 2) `content=${content##* ([:space:])}`

Linux

- Toepassing:
 - `content=${content##*+([[:digit:]])}`
 - String functie om te knippen
 - `#` een zo klein mogelijk string vinden
 - `##` een zo groot mogelijk string vinden
 - 1 of meer cijfers (equivalent met `[[:digit:]]+` in regex)
 - `*` = pattern = any string (including NULL)
 - Besluit:
 - Zoeken de grootst mogelijke string die start met eender wat en eindigt op 1 of meer getallen en knippen die weg.

Linux

- Toepassing:
 - `content=${content##*+([[:digit:]])}`
 - String functie om vooraan te knippen
 - `#` een zo klein mogelijk string vinden
 - `##` een zo groot mogelijk string vinden
 - 1 of meer cijfers (equivalent met `[[:digit:]]+` in regex)
 - `*` = pattern = any string (including NULL)
 - Besluit:
 - Zoeken de grootst mogelijke string die start met eender wat en eindigt op 1 of meer getallen en knippen die weg.

Linux

- Toepassing:
 - `content=${content##*([[:space:]])}`
 - String functie om vooraan te knippen
 - `#` een zo klein mogelijk string vinden
 - `##` een zo groot mogelijk string vinden
 - 0 of meer cijfers (zie extended patterns)
 - Besluit:
 - Zoeken naar 0 of meer spaties en deze wegnippen.

Hulpmiddelen

- <http://regexraptor.net/>
- <http://www.regexr.com/>
- CheatSheet in cursus...