

Regular Expressions Basics

re = regular expression.
greedy = always favour matching
lazy = always favour not matching

Any character

`.` (dot) Any character, except new line
`\.` Literal dot

Character classes (custom)

`[...]` Any character listed
`[^...]` Any character *not* listed
`[a-z]` Range, any character from a to z
`[-...]` To include dash (-), add it first

Boundaries

`^, $` [Start, End] of line
`\A, \Z, \z` [Start, End, End] of string
The difference between these occurs in multi-line mode.

`\b` Word boundary (Perl, Java, .Net)
`\B` Not word boundary
`\<, \>` [Start, End] of word[†] (Vim)
`\G` ... Match starts where the previous match ended
`\Qtext\E` Literal *text* inside regex

[†] Perl/Java/.Net:

`\<` → `(?<!\w) (?=\w)`
`\>` → `(?<=\w) (?!\w)`

(see *Lookaround* in next page).

Grouping and Backreferences

`(re)` Group *re* as a unit (Perl, Java, .Net)
`\(re\)` Vim without *very magic*
`\1, ..., \9` [1..9]-th matched ()

Substitutions

`$1, ..., $9` [1..9]-th matched ()
`$+` Last matched ()
`$&, &` Entire matched string (Perl, Vim)
`$`, $'` Text [before, after] match
`$_` Entire input text

Quantifiers (greedy)

(Perl/Java/.Net, Vim without *very magic*):

`re?, re\?` *re* optional (0 or 1)
`re*` Any quantity of *re* (0 or more)
`re+, re\+` At least one of *re* (1 or more)
`re{n}, re\{n}` *re* exactly *n* times
`re{n,}, re\{n,}` *re* at least *n* times
`re{n,m}, re\{n,m}` *re* at least *n* and at most *m* times

Quantifiers (lazy)

(Perl/Java/.Net, Vim without *very magic*):

`re?~, re\{-,1}` *re* optional (0 or 1)
`re*~, re\{-}` Any quantity of *re* (0 or more)
`re+~, re\{-1,}` At least one of *re* (1 or more)
`re{n}?, re\{-n}` *re* exactly *n* times
`re{n,}?, re\{-n,}` *re* at least *n* times
`re{n,m}?, re\{-n,m}` .. *re* between *n* and *m* times

Alternation (branching)

`re|re` *re* or *re* (Perl/Java/.Net)
`re\|re` *re* or *re* (Vim without *very magic*)
Note: In NFA engines alternatives are tried in order of appearance. So placing them in probability order prevents backtracking, improving efficiency.

Shorthands

`\s` White space (space, tab, etc.)[†]
`\S` Not `\s`
`\d` Digit (`[0-9]`)[†]
`\D` Not `\d`
`\w` Word char (`[a-zA-Z0-9_]`)[†]
`\W` Not `\w`
`\t, \n` Tab, New line
[†] Some flavours also recognise special Unicode characters in the same character group.
In Vim `_s` also includes new line.

Modes and Flags

`i` Ignore case[†] (Perl/Vim)
`m` Multi-line mode (Perl)
`s` Dot matches new line[‡] (Perl)
`x` Ignore white space[§] (Perl)
`g` Apply substitution in all occurrences (Perl/Vim)

[†] Java: `Pattern.CASE_INSENSITIVE`,
.Net: `RegexOptions.IgnoreCase`.
[‡] In Vim, use `_.` to make dot match new lines.
[§] Java: `Pattern.COMMENTS`,
.Net: `RegexOptions.IgnorePatternWhiteSpace`.

Regular Expressions Advanced

Lookaround

Lookahead = Check if *re* matches at current position, but do not consume.
Lookbehind = Check if *re* matches just before what follows, but do not consume.
Negative = Check if *re* does *not* match.

(Perl/Java/.Net, Vim):

<code>(?=re), re\@=</code>	Lookahead
<code>(?!re), re\@!</code>	Negative lookahead
<code>(?<=re), re\@<=</code>	Lookbehind
<code>(?<!re), re\@<!</code>	Negative lookbehind

Non-capturing parentheses and comments

Group as a unit, but do not capture for backreferencing.

<code>(?:re)</code>	Perl, Java, .Net
<code>\%(re\)</code>	Vim

<code>(?#free text)</code>	Comment inside regex
----------------------------	----------------------

Branch reset

`(?|re)` . Restart group numbering for each branch in *re*

Not available in Vim.

Named captures

<code>(?<name>re)</code>	Capture <i>re</i> under <i>name</i>
<code>\k<name></code>	Named backreference
<code>\${name}</code>	Named substitution (Perl)
<code>\$ {name}</code>	Named substitution (.Net)

Not available in Vim.

Atomic grouping

Match *re* without retry. In a NFA engine, discard all the possible backtracking states for the enclosed *re*.

<code>(?>re)</code>	Perl/Java/.Net
<code>re\@></code>	Vim

Possessive quantifiers

<code>re?+</code>	<code>(?>re?)</code>
<code>re*+</code>	<code>(?>re*)</code>
<code>re++</code>	<code>(?>re+)</code>
<code>re{n,m}+</code>	<code>(?>re{n,m})</code>

Not available in Vim (but can be done with atomic grouping).

End of previous match

`\G` Matches where the previous match ended
Useful when using the *g* mode.
In Perl, use `pos($str)` to know the current location where `\G` would match in `$str`.
Not available in Vim.

Conditional expressions

If *cond* then *re1*, else *re2*. *re2* is optional.

Perl/Java/.Net:

<code>(? (cond) re)</code>	Without <i>else re</i>
<code>(? (cond) re1 re2)</code>	With <i>else re</i>
<code>(? (n) ...)</code>	<i>cond</i> = <i>n</i> -th () matched

cond can be a *lookaround* expression (see the examples).
Not available in Vim.

Dynamic regex and embedded code

<code>(? {code})</code>	Match <i>re</i> built by <i>code</i> here
<code>{code}</code>	<i>code</i> that does anything

Useful for debugging regex by printing.

Not available in Vim.

Recursive expressions

<code>(?R)</code>	Repeat entire <i>re</i> here
<code>(?Rn)</code>	Repeat <i>re</i> captured under group <i>n</i> here

Not available in Vim.

Mode modifiers

(Perl/Java/.Net, Vim):

<code>(?i), \c</code>	Turn case-insensitive on
<code>(?-i), \C</code>	Turn case-insensitive off
<code>(?i:re)</code>	Be case insensitive for <i>re</i>

Regular Expressions and Characters

Unicode properties

Perl/Java/.Net:

<code>\p{L}</code>	Letters
<code>\p{M}</code>	Accent marks, etc.
<code>\p{Z}</code>	Spaces, etc.
<code>\p{S}</code>	Dingbats and symbols
<code>\p{N}</code>	Numeric characters
<code>\p{P}</code>	Punctuation characters
<code>\p{C}</code>	Everything else

Unicode sub-properties

<code>\p{Ll}</code>	Lower-case letters
<code>\p{Lu}</code>	Upper-case letters
<code>\p{Mn}</code>	Non-spacing mark (accents, ...)
<code>\p{Sm}</code>	Math symbol (−,+,÷,≤,...)
<code>\p{Sc}</code>	Currency symbol (¥,¢,€, \$,£,...)
<code>\p{Nd}</code>	Decimal digit
<code>\p{Nl}</code>	Letter number (mostly Roman numerals)
<code>\p{Pd}</code>	Dash punctuation
<code>\p{Ps}</code>	Open punctuation ([, {, ⟨, ...)
<code>\p{Pe}</code>	Close punctuation (], }, ⟩, ...)
<code>\p{Cc}</code>	.. ASCII and Latin-1 control characters (TAB, CR, LF, ...)

Unicode negated properties

<code>\P{...}</code>	Perl/Java/.Net
<code>\p{^...}</code>	Perl

Unicode blocks and scripts

<code>\p{InCyrillic}</code>	Cyrillic character (Perl, Java)
<code>\p{IsCyrillic}</code>	Cyrillic character (.Net)
<code>\p{Latin}</code>	Latin character
<code>\p{Greek}</code>	Greek character
<code>\p{Hebrew}</code>	Hebrew character

POSIX character classes

<code>[:alnum:]</code>	Alphanumeric characters
<code>[:alpha:]</code>	Alphabetic characters
<code>[:blank:]</code>	Space and tab
<code>[:cntrl:]</code>	Control characters
<code>[:digit:]</code>	Numeric characters
<code>[:graph:]</code>	Non-blank characters
<code>[:lower:]</code>	Lower-case alphabetic characters
<code>[:print:]</code>	... <code>[:graph:]</code> and the space character
<code>[:punct:]</code>	Punctuation characters
<code>[:space:]</code>	All whitespace characters
<code>[:upper:]</code>	Upper-case alphabetic characters
<code>[:xdigit:]</code>	Hexadecimal digits (<code>[0-9a-fA-F]</code>)

POSIX collating sequences

<code>[.span-ll.]</code>	Match “ll” as single character
<code>[.ch.]</code>	Match “ch” as single character
<code>[.eszet.]</code>	Match “ss” (“ß”) as single character

POSIX character equivalents

<code>[[=a=] [=n=]]</code>	Character <i>equivalents</i> for <i>a</i> and <i>n</i> (<i>a</i> , <i>á</i> , <i>à</i> , <i>ä</i> , ..., <i>n</i> , <i>ñ</i> , ...)
------------------------------	--

Special characters and other shorthands

<code>\char</code>	Literal <i>char</i>
<code>\t</code>	Tab (HT, TAB)
<code>\n</code>	New line (LF, NL)
<code>\r</code>	Carriage return (CR)
<code>\f</code>	Form feed (FF)
<code>\a</code>	Alarm (BEL)
<code>\e</code>	Escape (think troff) (ESC)
<code>\$\$</code>	Literal \$ (.Net)
<code>\l</code>	Lower-case next char
<code>\u</code>	Upper-case next char
<code>\Ltext\E</code>	Lower-case <i>text</i>
<code>\Utext\E</code>	Upper-case <i>text</i>
<code>\num</code>	Octal escape
<code>\xnum, \x{num}</code>	Hex escape
<code>\unum, \Unum</code>	Unicode escape
<code>\cx, \cX</code>	CTRL-X
<code>\N{U+263D}</code>	Unicode character
<code>\N{name}</code>	Character by Unicode <i>name</i>

Class set operations

<code>[[a-z]&&[^aeiou]]</code>	.. Class 1 except class 2 (Java)
<code>[[a-z]-[aeiou]]</code> Class 1 except class 2 (.Net)

Regular Expressions Efficiency

Notation: *worse* → *better*

General guidelines

- Avoid recompiling the *re*.
- Use non-capturing parentheses (but check they are indeed faster in your case!).
- Split into multiple *re*'s (or literal text search) if it can make the process faster.
- Use leading anchors (`^`, `\A`, `$`, `\Z`, etc.) to reduce the number of locations where the *re* is evaluated.
- Use atomic grouping and possessive quantifiers to avoid unnecessary backtracking (especially important in non-matching cases).
- Put the most-likely alternatives first (only affects traditional NFA engines).

Lazy quantifiers may be slower than greedy

Because they must jump between checking what the quantifier controls with checking what comes after. Example: `"(.*)"` to match a double-quoted string. Compare with the example in the next page.

Expose anchors

`^this|^that` → `^(?:this|that)`.

`abc$|123$` → `(?:abc|123)$` (only affects Perl).

re only tried at places where the anchor (i.e. `^`) matches.

Algebraic identities

Note: In this section all parenthesis are non-capturing. These are strictly theoretical equivalences.

`r|s = s|r` commutativity for alternation[†]
`r|(s|t) = (r|s)|t` associativity for alternation
`r|r = r` absorption of alternation
`r(st) = (rs)t` associativity for concatenation
`r(s|t) = rs|rt` left distributivity
`(s|t)r = sr|tr` right distributivity
`rε = r` identity for concatenation
`r*r* = r*` closure absorption
`r* = ε|r|rr|...` Kleene closure
`rr* = r*r`
`(r*)* = r*`
`(r*s*)* = (r*s*)*`
`(r*s*)* = (r|s)*`
`(rs)*r = r(sr)*`
`(r|s)* = (r*s)*r*`

[†] It is not the same in terms of performance for a NFA engine, where alternatives are tried in order of appearance.

Expose literal text

`this|that` → `th(?:is|at)`.

`-{5,7}` → `-----{0,2}`.

The engine can use a fast substring search (like the [Boyer-Moore](#) or the Hume-Sunday algorithms) to find the literal text (i.e. `th`) and then match the rest of the *re*.

Loop unrolling

Useful to optimise a regex with the following structure:

`(normal|special)*`. This can be changed to:

`opening normal*(?:special normal)* closing`, where `opening`, `normal`, `special` and `closing` are all regular expressions. It can be used to optimise any number of alternatives, like in `(re1|re2|...)*`.

Notes:

1. `normal` should be the most common case
2. the start of `normal` and `special` must not match at the same location.
3. `special` must not match the *empty string*.
4. `special` must be atomic.

See the examples.

Regular Expressions Examples 1

Match quoted string with escaped quotes

```
" ( [^\\"]++ | \\\. ) * +"
```

Note: Possessive quantifiers prevent very long execution times (due to the nested quantifiers `+` and `*`) when there is no matching.

Loop-unrolling version:

```
" [^\\"]*+ (?: \\. [^\\"]*+ )* +"
```

Conditional expressions

- Backreference as condition: Match a word optionally wrapped in `<...>`:

```
( < > ? \w+ ( ? {1} > )
```

- Lookaround as condition: Check for number only if prefixed with “NUM:”:

```
( ? ( ? <= NUM : ) \d+ | \w+ )
```

Branch reset

Always capture the number under `\1` (or `$1`):

```
( ? | Num : ( \d+ ) | Number : ( \d+ ) | N = ( \d+ ) )
```

Fix floating-point rounding problems

Use at most 3 decimal places. Change numbers like 3.27600000002828 into 3.276; or 4.120000000034 into 4.12:

```
s / ( \. \d \d [1-9] ?+ ) \d+ / $1 / g
```

Extract file name from path

Perl:

```
$path =~ m{ ([^/]*) $ };  
$file = $1;
```

Adding thousand separators to a number

Using lookahead:

```
s / ( ? <= \d ) ( ? = ( \d \d \d ) + ( ? ! \d ) ) / , / g
```

Without using lookahead (in Perl):

```
while ( $text =~ s / ( \d ) ( ( \d \d \d ) + \b ) / $1, $2 / g )  
{  
    # Just repeat until no match  
}
```

Match continuation line

Match a single “logical” line split into multiple lines by adding `\` at the end of each split. Example:

```
SRC=a.c b.c c.c \  
      d.c e.c
```

```
^ \w+ = ( [^ \n \\ ] | \\ . ) *
```

Loop-unrolling version:

```
^ \w+ =                # Leading field and '='  
(                     # Capture complete line  
    ( ? > [^ \n \\ ] * ) # normal  
    ( ? >              # special  
        [^ \n \\ ] *    # normal  
    ) *  
)
```

Match a date (month and day)

```
( ? |  
    ( Jan | Mar | May | Jul | Aug | Oct | Dec )  
    ( 31 | [123] 0 | [012] ? [1-9] )  
    |  
    ( Apr | Jun | Sep | Nov )  
    ( [123] 0 | [012] ? [1-9] )  
    |  
    Feb  
    ( [12] 0 | [012] ? [1-9] )  
)
```

Match time (am/pm)

```
( 1 [012] | [1-9] ) : [0-5] [0-9] ( am | pm )
```

Match time (24 hours)

```
[01] ? [0-9] | 2 [0-3]
```

or

```
[01] ? [4-9] | [012] ? [0-3]
```

Regular Expressions Examples 2

Parse a CSV file

This regex matches each field in a CSV line, supporting fields with or without surrounding double-quotes. In the former case, a double quote is represented by a pair of double quotes:

```
(?:^(,)|\s*
(?:|
    "
    (?: [^"] | "" )* ) # Slow
    "
|
    ( [^",]* )
)
```

The line with the “slow” comment can be made faster using loop unrolling, with normal = `[^"]` and special = `"`:

```
( (?:> [^"]* ) (?:> "" [^"]* )* )
```

Match IP address

Using Perl regex object for clarity:

```
$num = qr/[01]?[d\d]?[2[0-4]|d|25[0-5]]/;
(?:!\d)
($num)\.($num)\.($num)\.($num)
(?:!\d)
```

Match an email address

```
\b(
    \w[-.\w]* (?# user name)
    \@
    [-a-z0-9]+(\.[-a-z0-9]+)*
    \. (com|org|net|gov|edu|info)
    (\.[-a-z]{2})?
)\b
```

Match a URL

```
\<(
    (?:https?|ftp)://
    (?:i: [-a-z0-9] (?: [-a-z0-9]* [-a-z0-9] )? \.)+
    (?:|
        ((?-i: com|org|net|gov|edu|info)
        (?-i: [-a-z]{2})?)
        |
        ((?-i: [-a-z]{2}))
    )
    (:\d+)? (?# port number)
    \b
    (
        /
        [-a-z0-9_:@&?+=,./~*'%'$]*
        (?![.,?!]) (?# Not allowed at end)
    )?
)\>
```

Match closing XML tag

Using lazy quantifier: `((?!).)*?`

Using greedy quantifier: `((?!</?B>).)*`

Loop-unrolling with normal = `[^<]` and special = `(?!</?B><):`

```
<B>                # opening
(?:[<]*)          # any "normal"
(?:              # any amount of
    (?!</?B>)    # if not <B> or </B>
    <            # one "special"
    [<]*        # any normal
)*
</B>              # closing
```

Balanced parentheses (static)

Match up to n levels of nested parentheses:

```
# Using Perl regex objects
my $level0 = qr/\(((^()))*\)/x;
my $level1 = qr/\(((^())|level0)*\)/;
my $level2 = qr/\(((^())|level1)*\)/;
...
my $balpar = qr/\(((^())|levelN)*\)/;

# Using string to build regex
my $balpar = "\(((^())*)"; # level 0
foreach (1..$n) {
    $balpar = "\(((^())|$balpar)*\>";
}
}
```

Balanced parentheses (dynamic)

Match a function call with balanced parentheses:

```
my $balpar; # must be predefined
$balpar = qr/
    (?>
        [^()]+
        |
        \((??{ $balpar }) \)
    )*
/x;

if ($text =~ m/\b(\w+)\($balpar\)/) {
    print "function: $1, args: $2\n";
}
if (not $text =~ m/^ $balpar $/x) {
    print "mismatched parentheses\n";
}
```

Regular Expressions Usage in Programming Languages and Tools

Perl

```
# Highlight double words
# perl -w finddbl file.txt

# Chunk with dot-newline combination
$/ = ".\n"

while (<>) { # Put input "line" in $_
    next unless s{
        \b
        ( \w+ ) (?# grab word in $1 and \1 )
        (?# Any number of spaces or <TAGS> )
        (
            (?
                \s
                |
                <[>]+>
            )+
        )
        (\1\b) (?# match the first word again)
    }
    { \e[7m$1\e[m$2\e[7m$3\e[m] i g x;

    # Remove unmarked lines
    s/^(?![^\e]*\n)+//mg;

    # Insert file name
    s/^\$ARGV: /mg;

    # Print $_
    print;
}
```

Java

```
import java.util.regex.*;

String text;

// Extract subject
Pattern re = Pattern.compile(
    "^Subject: (.*)",
    Pattern.CASE_INSENSITIVE);
Matcher m = re.matcher(text);

if (m.find()) {
    subject = m.group(1);
}

// Insert prefix
Pattern re = Pattern.compile("^(.*)$");
re.matcher(text).replaceAll(">> $1");
```

Python

```
import re

r = re.compile("^Subject: (.*)",
               re.IGNORECASE)
m = r.search(text);
if m:
    subject = m.group(1)
```

.Net (C#)

```
using System.Text.RegularExpressions;

string text;

// Extract subject
Regex re = new Regex(
    "^Subject: (.*)",
    RegexOptions.IgnoreCase);
Match m = re.Match(text);

if (m.Success) {
    subject = m.Groups(1).Value;
}

// Insert prefix
Regex re = new Regex("^(.*)$");
re.Replace(text, ">> {1}");
```

Automated editing

```
sed -i .old -E 's/re/.../g' file
perl -p -i .old -e 's/re/.../g' file
```