

Ruby Regexp

a magical tool for text processing

Give me all leaves



Give me tree
with 2 branches



Sundeeep Agarwal

Table of contents

Preface	3
Prerequisites	3
Conventions	3
Acknowledgements	3
Feedback and Errata	4
Author info	4
License	4
Book version	4
Why is it needed?	5
How this book is organized	6
Regex introduction	7
Regex documentation	7
match? method	7
Regex literal reuse and interpolation	8
sub and gsub methods	9
Regex operators	9
Cheatsheet and Summary	10
Exercises	11
Anchors	14
String anchors	14
Line anchors	15
Word anchors	17
Cheatsheet and Summary	18
Exercises	19
Alternation and Grouping	21
OR conditional	21
Regex.union method	21
Grouping	22
Regex.source method	22
Precedence rules	23
Cheatsheet and Summary	24
Exercises	25
Escaping metacharacters	27
Escaping with \	27
Regex.escape method	27
Escaping delimiter	28
Escape sequences	28
Cheatsheet and Summary	29
Exercises	30

Preface

Scripting and automation tasks often need to extract particular portions of text from input data or modify them from one format to another. This book will help you learn Regular Expressions, a mini-programming language for all sorts of text processing needs.

The book heavily leans on examples to present features of regular expressions one by one. It is recommended that you manually type each example and experiment with them. Understanding both the nature of sample input string and the output produced is essential. As an analogy, consider learning to drive a bike or a car — no matter how much you read about them or listen to explanations, you need to practice a lot and infer your own conclusions. Should you feel that copy-paste is ideal for you, [code snippets are available chapter wise on GitHub](#).

Prerequisites

You should have prior experience working with Ruby, should know concepts like blocks, string formats, string methods, Enumerable, etc.

If you have prior experience with a programming language, but new to Ruby, check out my GitHub repository on [Ruby Scripting](#) before starting this book. That repository also includes a chapter on Regular Expressions which has been edited and expanded to create this book.

Conventions

- The examples presented here have been tested with **Ruby version 3.0.0** and includes features not available in earlier versions.
- Code snippets shown are copy pasted from `irb --simple-prompt` shell and modified for presentation purposes. Some commands are preceded by comments to provide context and explanations. Blank lines have been added to improve readability. `nil` return value is not shown for `puts` statements. Error messages are shortened. And so on.
- Unless otherwise noted, all examples and explanations are meant for **ASCII** characters.
- External links are provided for further reading throughout the book. Not necessary to immediately visit them. They have been chosen with care and would help, especially during re-reads.
- The [Ruby_Regexp repo](#) has all the code snippets and files used in examples and exercises and other details related to the book. If you are not familiar with `git` command, click the **Code** button on the webpage to get the files.

Acknowledgements

- [ruby-lang documentation](#) — manuals and tutorials
- [/r/ruby/](#) and [/r/regex/](#) — helpful forum for beginners and experienced programmers alike
- [stackoverflow](#) — for getting answers to pertinent questions on Ruby and regular expressions
- [tex.stackexchange](#) — for help on `pandoc` and `tex` related questions
- Cover image:
 - [draw.io](#)
 - [tree icon](#) by Gopi Doraisamy under [Creative Commons Attribution 3.0 Unported](#)
 - [wand icon](#) by [roundicons.com](#)
- [pngquant](#) and [svgcleaner](#) for optimizing images

- [Warning](#) and [Info](#) icons by [Amada44](#) under public domain
- [softwareengineering.stackexchange](#) and [skolakoda](#) for programming quotes
- [gmovchan](#) for spotting a typo

Special thanks to Allen Downey, an attempt at translating his book [Think Python](#) to [Think Ruby](#) gave me the confidence to publish my own book.

Feedback and Errata

I would highly appreciate if you'd let me know how you felt about this book, it would help to improve this book as well as my future attempts. Also, please do let me know if you spot any error or typo.

Issue Manager: https://github.com/learnbyexample/Ruby_Regexp/issues

E-mail: learnbyexample.net@gmail.com

Twitter: https://twitter.com/learn_byexample

Author info

Sundeeep Agarwal is a freelance trainer, author and mentor. His previous experience includes working as a Design Engineer at Analog Devices for more than 5 years. You can find his other works, primarily focused on Linux command line, text processing, scripting languages and curated lists, at <https://github.com/learnbyexample>. He has also been a technical reviewer for [Command Line Fundamentals](#) book and video course published by Packt.

List of books: <https://learnbyexample.github.io/books/>

License

This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License](#)

Code snippets are available under [MIT License](#)

Resources mentioned in Acknowledgements section above are available under original licenses.

Book version

2.6

See [Version_changes.md](#) to track changes across book versions.

Why is it needed?

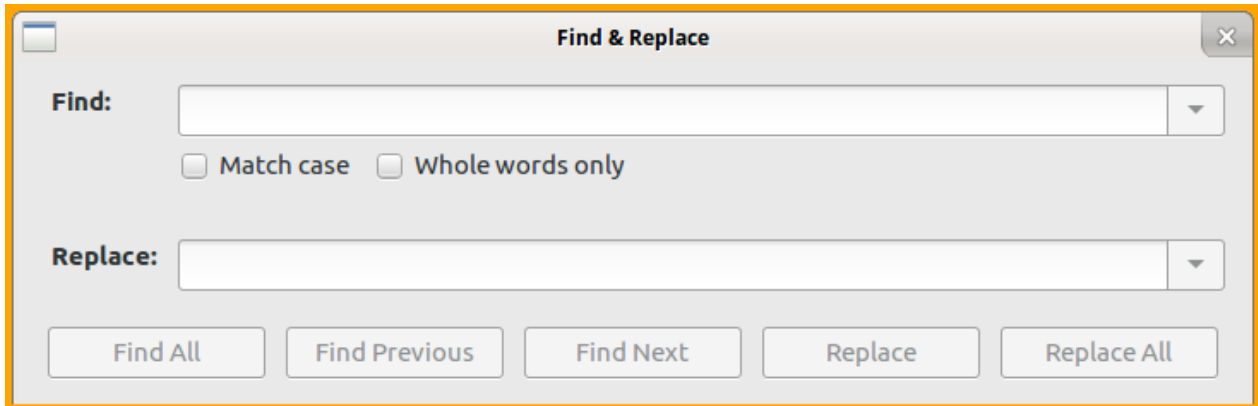
Regular Expressions is a versatile tool for text processing. You'll find them included as part of standard library of most programming languages that are used for scripting purposes. If not, you can usually find a third-party library. Syntax and features of regular expressions vary from language to language. Ruby's offering is based upon the [Onigmo regular expressions library](#).

The `String` class comes loaded with variety of methods to deal with text. So, what's so special about regular expressions and why would you need it? For learning and understanding purposes, one can view regular expressions as a mini programming language in itself, specialized for text processing. Parts of a regular expression can be saved for future use, analogous to variables and functions. There are ways to perform AND, OR, NOT conditionals. Operations similar to range and string repetition operators and so on.

Here's some common use cases:

- Sanitizing a string to ensure that it satisfies a known set of rules. For example, to check if a given string matches password rules.
- Filtering or extracting portions on an abstract level like alphabets, numbers, punctuation and so on.
- Qualified string replacement. For example, at the start or the end of a string, only whole words, based on surrounding text, etc.

You are likely to be familiar with graphical search and replace tool, like the screenshot shown below from LibreOffice Writer. **Match case**, **Whole words only**, **Replace** and **Replace All** are some of the basic features supported by regular expressions.



Another real world use case is password validation. The screenshot below is from GitHub sign up page. Performing multiple checks like **string length** and the **type of characters allowed** is another core feature of regular expressions.

A screenshot of the GitHub password validation form. It features a label 'Password *' followed by a text input field. Below the input field is a message: 'Make sure it's at least 15 characters OR at least 8 characters including a number and a lowercase letter. [Learn more.](#)'

Here's some articles on regular expressions to know about its history and the type of problems it is suited for.

- [The true power of regular expressions](#) — it also includes a nice explanation of what *regular* means in this context
- [softwareengineering: Is it a must for every programmer to learn regular expressions?](#)
- [softwareengineering: When you should NOT use Regular Expressions?](#)
- [codinghorror: Now You Have Two Problems](#)
- [wikipedia: Regular expression](#) — this article includes discussion on regular expressions as a formal language as well as details on various implementations

How this book is organized

The book introduces concepts one by one and exercises at the end of chapters will require only the features introduced until that chapter. Each concept is accompanied by multiple examples to cover various angles of usage and corner cases. As mentioned before, follow along the illustrations by typing out the code snippets manually. It is important to understand both the nature of the sample input string as well as the actual programming command used. There are two interlude chapters that give an overview of useful external resources and some more resources are collated in the final chapter.

- [Regexp introduction](#)
- [Anchors](#)
- [Alternation and Grouping](#)
- [Escaping metacharacters](#)
- [Dot metacharacter and Quantifiers](#)
- [Interlude: Tools for debugging and visualization](#)
- [Working with matched portions](#)
- [Character class](#)
- [Groupings and backreferences](#)
- [Interlude: Common tasks](#)
- [Lookarounds](#)
- [Modifiers](#)
- [Unicode](#)
- [Further Reading](#)

By the end of the book, you should be comfortable with both writing and reading regular expressions, how to debug them and know when to *avoid* them.

Regexp introduction

In this chapter, you'll get to know how to declare and use regexps. For some examples, the equivalent normal string method is shown for comparison. Regular expression features will be covered next chapter onwards. The main focus will be to get you comfortable with syntax and text processing examples. Three methods will be introduced in this chapter. The `match?` method to search if the input contains a string and the `sub` and `gsub` methods to substitute a portion of the input with something else.



This book will use the terms **regular expressions** and **regexp** interchangeably.

Regexp documentation

It is always a good idea to know where to find the documentation. Visit [ruby-doc: Regexp](#) for information on `Regexp` class, available methods, syntax, features, examples and more. Here's a quote:

Regular expressions (*regexps*) are patterns which describe the contents of a string. They're used for testing whether a string contains a given pattern, or extracting the portions that match. They are created with the `/pat/` and `%r{pat}` literals or the `Regexp.new` constructor.

match? method

First up, a simple example to test whether a string is part of another string or not. Normally, you'd use the `include?` method and pass a string as argument. For regular expressions, use the `match?` method and enclose the search string within `//` delimiters (regexp literal).

```
>> sentence = 'This is a sample string'

# check if 'sentence' contains the given string argument
>> sentence.include?('is')
=> true
>> sentence.include?('z')
=> false

# check if 'sentence' matches the pattern as described by the regexp argument
>> sentence.match?(/is/)
=> true
>> sentence.match?(/z/)
=> false
```

The `match?` method accepts an optional second argument which specifies the index to start searching from.

```
>> sentence = 'This is a sample string'

>> sentence.match?(/is/, 2)
```

```
=> true
>> sentence.match?(/is/, 6)
=> false
```

Some of the regular expressions functionality is enabled by passing modifiers, represented by an alphabet character. If you have used command line, modifiers are similar to command options, for example `grep -i` will perform case insensitive matching. It will be discussed in detail in [Modifiers](#) chapter. Here's an example for `i` modifier.

```
>> sentence = 'This is a sample string'

>> sentence.match?(/this/)
=> false
# 'i' is a modifier to enable case insensitive matching
>> sentence.match?(/this/i)
=> true
```

Regexp literal reuse and interpolation

The regexp literal can be saved in a variable. This helps to improve code clarity, pass around as method argument, enable reuse, etc.

```
>> pet = /dog/i
>> pet
=> /dog/i

>> 'They bought a Dog'.match?(pet)
=> true
>> 'A cat crossed their path'.match?(pet)
=> false
```

Similar to double quoted string literals, you can use interpolation and escape sequences in a regexp literal. See [ruby-doc: Strings](#) for syntax details on string escape sequences. Regexp literals have their own special escapes, which will be discussed in [Escape sequences](#) section.

```
>> "cat\tdog".match?(/\t/)
=> true
>> "cat\tdog".match?(/\a/)
=> false

>> greeting = 'hi'
>> /#{greeting} there/
=> /hi there/
>> /#{greeting.upcase} there/
=> /HI there/
>> /#{2**4} apples/
=> /16 apples/
```


sub and gsub methods

For search and replace, use `sub` or `gsub` methods. The `sub` method will replace only the first occurrence of the match, whereas `gsub` will replace all the occurrences. The regexp pattern to match against the input string has to be passed as the first argument. The second argument specifies the string which will replace the portions matched by the pattern.

```
>> greeting = 'Have a nice weekend'

# replace first occurrence of 'e' with 'E'
>> greeting.sub(/e/, 'E')
=> "HavE a nice weekend"
# replace all occurrences of 'e' with 'E'
>> greeting.gsub(/e/, 'E')
=> "HavE a nicE wEEkEnd"
```

Use `sub!` and `gsub!` methods for in-place substitution.

```
>> word = 'cater'

# this will return a string object, won't modify 'word' variable
>> word.sub(/cat/, 'wag')
=> "wager"
>> word
=> "cater"

# this will modify 'word' variable itself
>> word.sub!(/cat/, 'wag')
=> "wager"
>> word
=> "wager"
```

Regexp operators

Ruby also provides operators for regexp matching.

- `==` match operator returns index of the first match and `nil` if match is not found
- `!~` match operator returns `true` if string *doesn't* contain the given regexp and `false` otherwise
- `===` match operator returns `true` or `false` similar to the `match?` method

```
>> sentence = 'This is a sample string'

# can also use: /is/ == sentence
>> sentence == /is/
=> 2
>> sentence == /z/
=> nil

# can also use: /z/ !~ sentence
>> sentence !~ /z/
=> true
```

```
>> sentence !~ /is/
=> false
```

Just like `match?` method, both `==` and `!~` can be used in a conditional statement.

```
>> sentence = 'This is a sample string'

>> puts 'hi' if sentence =~ /is/
hi

>> puts 'oh' if sentence !~ /z/
oh
```

The `===` operator comes in handy with Enumerable methods like `grep` , `grep_v` , `all?` , `any?` , etc.

```
>> sentence = 'This is a sample string'

# regexp literal has to be on LHS and input string on RHS
>> /is/ === sentence
=> true
>> /z/ === sentence
=> false

>> words = %w[cat attempt tattle]
>> words.grep(/tt/)
=> ["attempt", "tattle"]
>> words.all?(/at/)
=> true
>> words.none?(/temp/)
=> false
```



A key difference from `match?` method is that these operators will also set regexp related [global variables](#).

Cheatsheet and Summary

Note	Description
ruby-doc: Regexp	Ruby Regexp documentation
Onigmo doc	Onigmo library documentation
<code>/pat/</code> or <code>%r{pat}</code>	regexp literal
<code>var = /pat/</code>	interpolation and escape sequences can also be used
<code>/pat1#{expr}pat2/</code>	save regexp literal in a variable
<code>s.match?(/pat/)</code>	use result of an expression to build regexp
<code>s.match?(/pat/, 3)</code>	check if string <code>s</code> matches the pattern <code>/pat/</code>
<code>/pat/i</code>	returns <code>true</code> or <code>false</code>
	optional 2nd argument changes starting index of search
	modifier <code>i</code> matches alphabets case insensitively

Note	Description
<code>s.sub(/pat/, 'replace')</code>	search and replace first matching occurrence use <code>gsub</code> to replace all occurrences use <code>sub!</code> and <code>gsub!</code> for in-place substitution
<code>s =~ /pat/</code> or <code>/pat/ =~ s</code>	returns index of first match or <code>nil</code>
<code>s !~ /pat/</code> or <code>/pat/ !~ s</code>	returns <code>true</code> if no match or <code>false</code>
<code>/pat/ === s</code>	returns <code>true</code> or <code>false</code> similar to <code>match?</code> these operators will also set regexp global variables

This chapter introduced the `Regexp` class and methods `match?`, `sub` and `gsub` were discussed. You also learnt how to save and reuse regexp literals, how to specify modifiers and how to use regexp operators.

You might wonder why there are so many ways to test matching condition with regexps. The most common approach is to use `match?` method in a conditional statement. If you need position of match, use `=~` operator or `index` method. The `===` operator is usually relevant in Enumerable methods. Usage of global variables will be covered in later chapters. The `=~` and `!~` operators are also prevalent in command line usage (see my [Ruby one liners](#) tutorial for examples).

The next section has exercises to test your understanding of the concepts introduced in this chapter. Please do solve them before moving on to the next chapter.

Exercises



Refer to [exercises folder](#) for input files required to solve the exercises.



All the exercises are also collated together in one place at [Exercises.md](#). For solutions, see [Exercise_solutions.md](#).

a) Check whether the given strings contain `0xB0`. Display a boolean result as shown below.

```
>> line1 = 'start address: 0xA0, func1 address: 0xC0'
>> line2 = 'end address: 0xFF, func2 address: 0xB0'

>> line1.match?()      ##### add your solution here
=> false
>> line2.match?()      ##### add your solution here
=> true
```

b) For the given input file, print all lines containing the string `two`.

```
# note that expected output shown here is wrapped to fit pdf width
>> filename = 'programming_quotes.txt'

>> word = ##### add your solution here
```

```
>> puts File.foreach(filename).grep(word)
"Some people, when confronted with a problem, think - I know, I'll use regular
expressions. Now they have two problems" by Jamie Zawinski
"So much complexity in software comes from trying to make one thing do two
things" by Ryan Singer
```

c) Replace all occurrences of `5` with `five` for the given string.

```
>> ip = 'They ate 5 apples and 5 oranges'

>> ip.gsub(/, 'five')      ##### add your solution here
=> "They ate five apples and five oranges"
```

d) Replace first occurrence of `5` with `five` for the given string.

```
>> ip = 'They ate 5 apples and 5 oranges'

>> ip.sub(/, 'five')      ##### add your solution here
=> "They ate five apples and 5 oranges"
```

e) For the given array, filter all elements that do *not* contain `e` .

```
>> items = %w[goal new user sit eat dinner]

>> items.grep_v(/)      ##### add your solution here
=> ["goal", "sit"]
```

f) Replace all occurrences of `note` irrespective of case with `X` .

```
>> ip = 'This note should not be NoTeD'

>> ip.gsub(/, 'X')      ##### add your solution here
=> "This X should not be XD"
```

g) For the given input string, print all lines NOT containing the string `2`

```
'> purchases = %q{items qty
'> apple 24
'> mango 50
'> guava 42
'> onion 31
>> water 10}

>> num = //      ##### add your solution here

>> puts purchases.each_line.grep_v(num)
items qty
mango 50
onion 31
water 10
```

h) For the given array, filter all elements that contains either `a` or `w` .

```
>> items = %w[goal new user sit eat dinner]

>> items.filter { }      ##### add your solution here
=> ["goal", "new", "eat"]
```

i) For the given array, filter all elements that contains both `e` and `n` .

```
>> items = %w[goal new user sit eat dinner]

>> items.filter { }      ##### add your solution here
=> ["new", "dinner"]
```

j) For the given string, replace `0xA0` with `0x7F` and `0xC0` with `0x1F` .

```
>> ip = 'start address: 0xA0, func1 address: 0xC0'

##### add your solution here
=> "start address: 0x7F, func1 address: 0x1F"
```

k) Find the starting index of the first occurrence of `is` for the given input string.

```
>> ip = 'match this after the history lesson'

##### add your solution here
=> 8
```

Anchors

Now that you're familiar with regexp syntax and some of the methods, the next step is to know about the special features of regular expressions. In this chapter, you'll be learning about qualifying a pattern. Instead of matching anywhere in the given input string, restrictions can be specified. For now, you'll see the ones that are already part of regular expression features. In later chapters, you'll learn how to define your own rules for restriction.

These restrictions are made possible by assigning special meaning to certain characters and escape sequences. The characters with special meaning are known as **metacharacters** in regexp parlance. In case you need to match those characters literally, you need to escape them with a `\` character (discussed in [Escaping metacharacters](#) section).

String anchors

This restriction is about qualifying a regexp to match only at start or end of an input string. These provide functionality similar to the string methods `start_with?` and `end_with?`. There are three different escape sequences related to string level regexp anchors. First up is `\A` which restricts the matching to the start of string.

```
# \A is placed as a prefix to the search term
>> 'cater'.match?(/\Acat/)
=> true
>> 'concatenation'.match?(/\Acat/)
=> false

>> "hi hello\ntop spot".match?(/\Ahi/)
=> true
>> "hi hello\ntop spot".match?(/\Atop/)
=> false
```

To restrict the match to the end of string, `\z` is used.

```
# \z is placed as a suffix to the search term
>> 'spare'.match?(/are\z/)
=> true
>> 'nearest'.match?(/are\z/)
=> false

>> words = %w[surrender unicorn newer door empty eel pest]
>> words.grep(/er\z/)
=> ["surrender", "newer"]
>> words.grep(/t\z/)
=> ["pest"]
```

There is another end of string anchor `\Z`. It is similar to `\z` but if newline is the last character, then `\Z` allows matching just before the newline character.

```
# same result for both \z and \Z
# as there is no newline character at the end of string
>> 'dare'.sub(/are\z/, 'X')
=> "dX"
```

```
>> 'dare'.sub(/are\Z/, 'X')
=> "dX"

# different results as there is a newline character at the end of string
>> "dare\n".sub(/are\z/, 'X')
=> "dare\n"
>> "dare\n".sub(/are\Z/, 'X')
=> "dX\n"
```

Combining both the start and end string anchors, you can restrict the matching to the whole string. Similar to comparing strings using the `==` operator.

```
>> 'cat'.match?(/\Acat\z/)
=> true
>> 'cater'.match?(/\Acat\z/)
=> false
>> 'concatenation'.match?(/\Acat\z/)
=> false
```

The anchors can be used by themselves as a pattern. Helps to insert text at the start or end of string, emulating string concatenation operations. These might not feel like useful capability, but combined with other regexp features they become quite a handy tool.

```
>> 'live'.sub(/\A/, 're')
=> "relive"
>> 'send'.sub(/\A/, 're')
=> "resend"

>> 'cat'.sub(/\z/, 'er')
=> "cater"
>> 'hack'.sub(/\z/, 'er')
=> "hacker"
```

Line anchors

A string input may contain single or multiple lines. The newline character `\n` is used as the line separator. There are two line anchors, `^` metacharacter for matching the start of line and `$` for matching the end of line. If there are no newline characters in the input string, these will behave same as the `\A` and `\z` anchors respectively.

```
>> pets = 'cat and dog'

>> pets.match?(/^cat/)
=> true
>> pets.match?(/^dog/)
=> false

>> pets.match?(/dog$/)
=> true

>> pets.match?(/^dog$/)
```

```
=> false
```

Here's some multiline examples to distinguish line anchors from string anchors.

```
# check if any line in the string starts with 'top'
>> "hi hello\ntop spot".match?(/^top/)
=> true

# check if any line in the string ends with 'er'
>> "spare\ndare".match?(/er$/)
=> false

# filter all lines ending with 'are'
>> "spare\ndare".each_line.grep(/are$/)
=> ["spare\n", "dare"]

# check if any complete line in the string is 'par'
>> "spare\ndare".match?(/^par$/)
=> true
```

Just like string anchors, you can use the line anchors by themselves as a pattern. `gsub` and `puts` will be used here to better illustrate the transformation. The `gsub` method returns an Enumerator if you don't specify a replacement string nor pass a block. That paves way to use all those wonderful Enumerator and Enumerable methods.

```
>> str = "catapults\nconcatenate\ncat"

>> puts str.gsub(/^/, '1: ')
1: catapults
1: concatenate
1: cat
>> puts str.gsub(/^/).with_index(1) { |m, i| "#{i}: " }
1: catapults
2: concatenate
3: cat

>> puts str.gsub(/$/, '.')
catapults.
concatenate.
cat.
```

If there is a newline character at the end of string, there is an additional end of line match but no additional start of line match.

```
>> puts "1\n2\n".gsub(/^/, 'foo ')
foo 1
foo 2
>> puts "1\n\n".gsub(/^/, 'foo ')
foo 1
foo
```



```
# note the number of lines in output
>> puts "1\n2\n".gsub(/$/, ' baz')
1 baz
2 baz
  baz
>> puts "1\n\n".gsub(/$/, ' baz')
1 baz
  baz
  baz
```



If you are dealing with Windows OS based text files, you'll have to convert `\r\n` line endings to `\n` first. Which is easily handled by many of the Ruby methods. For example, you can specify which line ending to use for `File.open` method, the `split` string method handles all whitespaces by default and so on. Or, you can handle `\r` as optional character with quantifiers (see [Greedy quantifiers](#) section).

Word anchors

The third type of restriction is word anchors. Alphabets (irrespective of case), digits and the underscore character qualify as word characters. You might wonder why there are digits and underscores as well, why not only alphabets? This comes from variable and function naming conventions — typically alphabets, digits and underscores are allowed. So, the definition is more oriented to programming languages than natural ones.

The escape sequence `\b` denotes a word boundary. This works for both the start of word and end of word anchoring. Start of word means either the character prior to the word is a non-word character or there is no character (start of string). Similarly, end of word means the character after the word is a non-word character or no character (end of string). This implies that you cannot have word boundary `\b` without a word character.

```
>> words = 'par spar apparent spare part'

# replace 'par' irrespective of where it occurs
>> words.gsub(/par/, 'X')
=> "X sX apXent sXe Xt"
# replace 'par' only at the start of word
>> words.gsub(/\bpar/, 'X')
=> "X spar apparent spare Xt"
# replace 'par' only at the end of word
>> words.gsub(/par\b/, 'X')
=> "X sX apparent spare part"
# replace 'par' only if it is not part of another word
>> words.gsub(/\bpar\b/, 'X')
=> "X spar apparent spare part"
```

You can get lot more creative with using word boundary as a pattern by itself:

```
# space separated words to double quoted csv
# note the use of 'tr' string method
```

```
>> puts words.gsub(/\b/, '').tr(' ', ',')
"par","spar","apparent","spare","part"

>> '-----hello-----'.gsub(/\b/, ' ')
=> "----- hello -----"

# make a programming statement more readable
# shown for illustration purpose only, won't work for all cases
>> 'foo_baz=num1+35*42/num2'.gsub(/\b/, ' ')
=> " foo_baz = num1 + 35 * 42 / num2 "

# excess space at start/end of string can be stripped off
# later you'll learn how to add a qualifier so that strip is not needed
>> 'foo_baz=num1+35*42/num2'.gsub(/\b/, ' ').strip
=> "foo_baz = num1 + 35 * 42 / num2"
```

The word boundary has an opposite anchor too. `\B` matches wherever `\b` doesn't match. This duality will be seen with some other escape sequences too. Negative logic is handy in many text processing situations. But use it with care, you might end up matching things you didn't intend!

```
>> words = 'par spar apparent spare part'

# replace 'par' if it is not start of word
>> words.gsub(/\Bpar/, 'X')
=> "par sX apXent sXe part"

# replace 'par' at the end of word but not whole word 'par'
>> words.gsub(/\Bpar\b/, 'X')
=> "par sX apparent spare part"

# replace 'par' if it is not end of word
>> words.gsub(/par\b/, 'X')
=> "par spar apXent sXe Xt"

# replace 'par' if it is surrounded by word characters
>> words.gsub(/\Bpar\b/, 'X')
=> "par spar apXent sXe part"
```

Here's some standalone pattern usage to compare and contrast the two word anchors.

```
>> 'copper'.gsub(/\b/, ':')
=> ":copper:"

>> 'copper'.gsub(/\B/, ':')
=> "c:o:p:p:e:r"

>> '-----hello-----'.gsub(/\b/, ' ')
=> "----- hello -----"

>> '-----hello-----'.gsub(/\B/, ' ')
=> " - - - - -h e l l o- - - - -"
```

Cheatsheet and Summary

Note	Description
<code>\A</code>	restricts the match to the start of string
<code>\z</code>	restricts the match to the end of string
<code>\Z</code>	restricts the match to end or just before newline at the end of string
<code>\n</code>	line separator, dos-style files need special attention
metacharacter	characters with special meaning in regexp
<code>^</code>	restricts the match to the start of line
<code>\$</code>	restricts the match to the end of line
<code>\b</code>	restricts the match to the start/end of words word characters: alphabets, digits, underscore
<code>\B</code>	matches wherever <code>\b</code> doesn't match

In this chapter, you've begun to see building blocks of regular expressions and how they can be used in interesting ways. At the same time, regular expression is but another tool for text processing. Often, you'd get simpler solution by combining regular expressions with other string and Enumerable methods. Practice, experience and imagination would help you construct creative solutions. In coming chapters, you'll see more applications of anchors as well as `\G` [anchor](#) which is best understood in combination with other regexp features.

Exercises

a) Check if the given strings start with `be` .

```
>> line1 = 'be nice'
>> line2 = '"best!'"
>> line3 = 'better?'
>> line4 = 'oh no\nbear spotted'

>> pat = ##### add your solution here

>> pat.match?(line1)
=> true
>> pat.match?(line2)
=> false
>> pat.match?(line3)
=> true
>> pat.match?(line4)
=> false
```

b) For the given input string, change only whole word `red` to `brown`

```
>> words = 'bred red spread credible'

>> words.gsub() ##### add your solution here
=> "bred brown spread credible"
```

c) For the given input array, filter all elements containing `42` surrounded by word characters.

```
>> items = ['hi42bye', 'nice1423', 'bad42', 'cool_42a', 'fake4b']
```

```
>> items.grep()      ##### add your solution here
=> ["hi42bye", "nice1423", "cool_42a"]
```

d) For the given input array, filter all elements that start with `den` or end with `ly`

```
>> items = ['lovely', "1\ndentist", '2 lonely', 'eden', "fly\n", 'dent']

>> items.filter { }      ##### add your solution here
=> ["lovely", "2 lonely", "dent"]
```

e) For the given input string, change whole word `mall` to `1234` only if it is at start of line.

```
'> para = %q{ball fall wall tall
'> mall call ball pall
'> wall mall ball fall
>> {mallet wallet malls}

>> puts para.gsub()      ##### add your solution here
ball fall wall tall
1234 call ball pall
wall mallball fall
mallet wallet malls
```

f) For the given array, filter all elements having a line starting with `den` or ending with `ly`.

```
>> items = ['lovely', "1\ndentist", '2 lonely', 'eden', "fly\nfar", 'dent']

>> items.filter { }      ##### add your solution here
=> ["lovely", "1\ndentist", "2 lonely", "fly\nfar", "dent"]
```

g) For the given input array, filter all whole elements `12\nthree` irrespective of case.

```
>> items = ["12\nthree\n", "12\nThree", "12\nthree\n4", "12\nthree"]

>> items.grep()      ##### add your solution here
=> ["12\nThree", "12\nthree"]
```

h) For the given input array, replace `hand` with `X` for all words that start with `hand` followed by at least one word character.

```
>> items = %w[handed hand handy unhanded handle hand-2]

>> items.map { }      ##### add your solution here
=> ["Xed", "hand", "Xy", "unhanded", "Xle", "hand-2"]
```

i) For the given input array, filter all elements starting with `h`. Additionally, replace `e` with `X` for these filtered elements.

```
>> items = %w[handed hand handy unhanded handle hand-2]

>> items.filter_map { }      ##### add your solution here
=> ["handXd", "hand", "handy", "handlX", "hand-2"]
```

Alternation and Grouping

Many a times, you want to check if the input string matches multiple patterns. For example, whether an object's color is *green* or *blue* or *red*. In programming terms, you need to perform OR conditional. This chapter will show how to use alternation for such cases. These patterns can have some common elements between them, in which case grouping helps to form terser expressions. This chapter will also discuss the precedence rules used to determine which alternation wins.

OR conditional

A conditional expression combined with logical OR evaluates to `true` if any of the condition is satisfied. Similarly, in regular expressions, you can use `|` metacharacter to combine multiple patterns to indicate logical OR. The matching will succeed if any of the alternate pattern is found in the input string. These alternatives have the full power of a regular expression, for example they can have their own independent anchors. Here's some examples.

```
# match either 'cat' or 'dog'
>> 'I like cats'.match?(/cat|dog/)
=> true
>> 'I like dogs'.match?(/cat|dog/)
=> true
>> 'I like parrots'.match?(/cat|dog/)
=> false

# replace either 'cat' at start of string or 'cat' at end of word
>> 'catapults concatenate cat scat'.gsub(/\Acat|cat\b/, 'X')
=> "Xapults concatenate X sX"

# replace either 'cat' or 'dog' or 'fox' with 'mammal'
>> 'cat dog bee parrot fox'.gsub(/cat|dog|fox/, 'mammal')
=> "mammal mammal bee parrot mammal"
```

Regexp.union method

You might infer from above examples that there can be cases where lots of alternation is required. The `Regexp.union` method can be used to build the alternation list automatically. It accepts an array as argument or a list of comma separated arguments.

```
>> Regexp.union('car', 'jeep')
=> /car|jeep/

>> words = %w[cat dog fox]
>> pat = Regexp.union(words)
>> pat
=> /cat|dog|fox/
>> 'cat dog bee parrot fox'.gsub(pat, 'mammal')
=> "mammal mammal bee parrot mammal"
```



In the above examples, the elements do not contain any special regular expression characters. Strings having metacharacters will be discussed in [Regexp.escape method](#) section.

Grouping

Often, there are some common things among the regexp alternatives. It could be common characters or regexp qualifiers like the anchors. In such cases, you can group them using a pair of parentheses metacharacters. Similar to $a(b+c)d = abd+acd$ in maths, you get $a(b|c)d = abd|acd$ in regular expressions.

```
# without grouping
>> 'red reform read arrest'.gsub(/reform|rest/, 'X')
=> "red X read arX"

# with grouping
>> 'red reform read arrest'.gsub(/re(form|st)/, 'X')
=> "red X read arX"

# without grouping
>> 'par spare part party'.gsub(/\bpar\b|\bpart\b/, 'X')
=> "X spare X party"

# taking out common anchors
>> 'par spare part party'.gsub(/\b(par|part)\b/, 'X')
=> "X spare X party"

# taking out common characters as well
# you'll later learn a better technique instead of using empty alternate
>> 'par spare part party'.gsub(/\bpar(|t)\b/, 'X')
=> "X spare X party"
```



There's plenty more features to grouping than just forming terser regexp. It will be discussed as they become relevant in coming chapters.

Regexp.source method

The `Regexp.source` method helps to interpolate a regexp literal inside another regexp. For example, adding anchors to alternation list created using the `Regexp.union` method.

```
>> words = %w[cat par]
>> alt = Regexp.union(words)
>> alt
=> /cat|par/
>> alt_w = /\b(#{alt.source})\b/
>> alt_w
=> /\b(cat|par)\b/

>> 'cater cat concatenate par spare'.gsub(alt, 'X')
=> "Xer X conXenate X sXe"
```

```
>> 'cater cat concatenate par spare'.gsub(alt_w, 'X')
=> "cater X concatenate X spare"
```



The above example will work without `Regexp.source` method too, but you'll see that `/\b(#{alt})\b/` gives `/\b((?-mix:cat|par))\b/` instead of `/\b(cat|par)\b/`. Their meaning will be explained in [Modifiers](#) chapter.

Precedence rules

There's some tricky situations when using alternation. If it is used for testing a match to get `true/false` against a string input, there is no ambiguity. However, for other things like string replacement, it depends on a few factors. Say, you want to replace either `are` or `spared` — which one should get precedence? The bigger word `spared` or the substring `are` inside it or based on something else?

In Ruby, the regexp alternative which matches earliest in the input string gets precedence. Regexp operator `==` helps to illustrate this concept.

```
>> words = 'lion elephant are rope not'

>> words =~ /on/
=> 2
>> words =~ /ant/
=> 10

# starting index of 'on' < index of 'ant' for given string input
# so 'on' will be replaced irrespective of order of regexp
>> words.sub(/on|ant/, 'X')
=> "liX elephant are rope not"
>> words.sub(/ant|on/, 'X')
=> "liX elephant are rope not"
```

So, what happens if two or more alternatives match on same index? The precedence is then left to right in the order of declaration.

```
>> mood = 'best years'

>> mood =~ /year/
=> 5
>> mood =~ /years/
=> 5

# starting index for 'year' and 'years' will always be same
# so, which one gets replaced depends on the order of alternation
>> mood.sub(/year|years/, 'X')
=> "best Xs"
>> mood.sub(/years|year/, 'X')
=> "best X"
```

Another example with `gsub` to drive home the issue.

```
>> words = 'ear xerox at mare part learn eye'

# this is going to be same as: gsub(/ar/, 'X')
>> words.gsub(/ar|are|art/, 'X')
=> "eX xerox at mXe pXt leXn eye"

# this is going to be same as: gsub(/are|ar/, 'X')
>> words.gsub(/are|ar|art/, 'X')
=> "eX xerox at mX pXt leXn eye"

# phew, finally this one works as needed
>> words.gsub(/are|art|ar/, 'X')
=> "eX xerox at mX pX leXn eye"
```

If you do not want substrings to sabotage your replacements, a robust workaround is to sort the alternations based on length, longest first.

```
>> words = %w[hand handy handfu]

>> alt = Regexp.union(words.sort_by { |w| -w.length })
>> alt
=> /handful|handy|hand/

>> 'hands handfu handed handy'.gsub(alt, 'X')
=> "Xs X Xed X"

# without sorting, order will come into play
>> 'hands handfu handed handy'.gsub(Regexp.union(words), 'X')
=> "Xs Xful Xed Xy"
```

Cheatsheet and Summary

Note	Description
	multiple regexp combined as conditional OR
Regexp.union(array)	each alternative can have independent anchors
()	programmatically combine multiple strings/regexps
a(b c)d	group pattern(s)
	same as abd acd
/#{pat.source}/	interpolate a regexp literal inside another regexp
Alternation precedence	pattern which matches earliest in the input gets precedence
	tie-breaker is left to right if patterns have same starting location
	robust solution: sort the alternations based on length, longest first
	for ex: Regexp.union(words.sort_by { w -w.length })

So, this chapter was about specifying one or more alternate matches within the same regexp using `|` metacharacter. Which can further be simplified using `()` grouping if there are

common aspects among the alternations. Among the alternations, earliest matching pattern gets precedence. Left to right ordering is used as a tie-breaker if multiple alternations match starting from the same location. You also learnt couple of `Regexp` methods that help to programmatically construct a regexp literal.

Exercises

a) For the given input array, filter all elements that start with `den` or end with `ly`

```
>> items = ['lovely', "1\ndentist", '2 lonely', 'eden', "fly\n", 'dent']

>> items.grep()      ##### add your solution here
=> ["lovely", "2 lonely", "dent"]
```

b) For the given array, filter all elements having a line starting with `den` or ending with `ly`

```
>> items = ['lovely', "1\ndentist", '2 lonely', 'eden', "fly\nfar", 'dent']

>> items.grep()      ##### add your solution here
=> ["lovely", "1\ndentist", "2 lonely", "fly\nfar", "dent"]
```

c) For the given input strings, replace all occurrences of `removed` or `reed` or `received` or `refused` with `X`.

```
>> s1 = 'creed refuse removed read'
>> s2 = 'refused reed redo received'

>> pat =      ##### add your solution here

>> s1.gsub(pat, 'X')
=> "cX refuse X read"
>> s2.gsub(pat, 'X')
=> "X X redo X"
```

d) For the given input strings, replace all matches from the array `words` with `A`.

```
>> s1 = 'plate full of slate'
>> s2 = "slated for later, don't be late"
>> words = %w[late later slated]

>> pat =      ##### add your solution here

>> s1.gsub(pat, 'A')
=> "pA full of sA"
>> s2.gsub(pat, 'A')
=> "A for A, don't be A"
```

e) Filter all whole elements from the input array `items` that exactly matches any of the elements present in the array `words`.

```
>> items = ['slate', 'later', 'plate', 'late', 'slates', 'slated']
>> words = %w[late later slated]
```

```
>> pat = ##### add your solution here

>> items.grep(pat)
=> ["later", "late"]
```

Escaping metacharacters

This chapter will show how to match metacharacters literally, for manually as well as programmatically constructed patterns. You'll also learn about escape sequences supported by regexp and how they differ from strings.

Escaping with \

You have seen a few metacharacters and escape sequences that help to compose a regexp literal. To match the metacharacters literally, i.e. to remove their special meaning, prefix those characters with a `\` character. To indicate a literal `\` character, use `\\`.

To spice up the examples a bit, block form has been used below to modify the matched portion of string with an expression. In later chapters, you'll see more ways to work directly with matched portions.

```
# even though ^ is not being used as anchor, it won't be matched literally
>> 'a^2 + b^2 - C*3'.match?(/b^2/)
=> false

# escaping will work
>> 'a^2 + b^2 - C*3'.gsub(/(a|b)\^2/) { |m| m.upcase }
=> "A^2 + B^2 - C*3"

# match ( or ) literally
>> '(a*b) + c'.gsub(/(\(|\|)/, '')
=> "a*b + c"

>> '\learn\by\example'.gsub(/\\/, '/')
=> "/learn/by/example"
```

As emphasized earlier, regular expression is just another tool to process text. Some examples and exercises presented in this book can be solved using normal string methods as well. For real world use cases, ask yourself first if regexp is needed at all?

```
>> eqn = 'f*(a^b) - 3*(a^b)'

# straightforward search and replace, no need regexp shenanigans
>> eqn.gsub('(a^b)', 'c')
=> "f*c - 3*c"
```

Regexp.escape method

How to escape all the metacharacters when a regexp is constructed dynamically? Relax, `Regexp.escape` method has got you covered. No need to manually take care of all the metacharacters or worry about changes in future versions.

```
>> eqn = 'f*(a^b) - 3*(a^b)'
>> expr = '(a^b)'

>> puts Regexp.escape(expr)
\ (a\^b\)
```

```
# replace only at the end of string
>> eqn.sub(/#{Regexp.escape(expr)}\z/, 'c')
=> "f*(a^b) - 3*c"
```

The `Regexp.union` method automatically applies escaping for string arguments.

```
# array of strings, assume alternation precedence sorting isn't needed
>> terms = %w[a_42 (a^b) 2|3]

>> pat = Regexp.union(terms)
>> pat
=> /a_42|\(a^b\)|2|3/

>> 'ba_423 (a^b)c 2|3 a^b'.gsub(pat, 'X')
=> "bX3 Xc X a^b"
```

`Regexp.union` will also take care of mixing string and regexp patterns correctly. `(?-mix:` seen in the output below will be explained in the [Modifiers](#) chapter.

```
>> Regexp.union(/^cat|dog$/, 'a^b')
=> /(?-mix:^cat|dog$)|a^b/
```

Escaping delimiter

Another character to keep track for escaping is the delimiter used to define the regexp literal. Or, you can use a different delimiter than `/` to define a regexp literal using `%r` to avoid escaping. Also, you need not worry about unescaped delimiter inside `#{} interpolation`.

```
>> path = '/abc/123/foo/baz/ip.txt'

# \/ is also known as 'leaning toothpick syndrome'
>> path.sub(/\A/abc\123\\/, '~/')
=> "~/foo/baz/ip.txt"

# a different delimiter improves readability and reduces typos
>> path.sub(%r#\A/abc/123/#, '~/')
=> "~/foo/baz/ip.txt"
```

Escape sequences

In regexp literals, characters like tab and newline can be expressed using escape sequences as `\t` and `\n` respectively. These are similar to how they are treated in normal string literals (see [ruby-doc: Strings](#) for details). However, escapes like `\b` (word boundary) and `\s` (see [Escape sequence character sets](#) section) are different for regexps. And octal escapes `\nnn` have to be three digits to avoid conflict with [Backreferences](#).

```
>> "a\tb\tc".gsub(/\t/, ':')
=> "a:b:c"

>> "1\n2\n3".gsub(/\n/, ' ')
=> "1 2 3"
```

If an escape sequence is not defined, it'll match the character that is escaped. For example, `\%` will match `%` and not `\` followed by `%`.

```
>> 'h%x'.match?(/h%x/)
=> true
>> 'h\x'.match?(/h%x/)
=> false

>> 'hello'.match?(/\l/)
=> true
```

If you represent a metacharacter using escapes, it will be treated literally instead of its metacharacter feature.

```
# \x20 is hexadecimal for space character
>> 'h e l l o'.gsub(/\x20/, '')
=> "hello"

# \053 is octal for + character
>> 'a+b'.match?(/a\053b/)
=> true

# \x7c is '|' character
>> '12|30'.gsub(/2\x7c3/, '5')
=> "150"
>> '12|30'.gsub(/2|3/, '5')
=> "15|50"
```



See [ASCII code table](#) for a handy cheatsheet with all the ASCII characters and their hexadecimal representation.

[Codepoints and Unicode escapes](#) section will discuss escapes for unicode characters.

Cheatsheet and Summary

Note	Description
<code>\</code>	prefix metacharacters with <code>\</code> to match them literally
<code>\\</code>	to match <code>\</code> literally
<code>Regexp.escape(s)</code>	automatically escape all metacharacters for string <code>s</code>
<code>Regexp.union</code>	also automatically escapes string arguments
<code>%r</code>	helps to avoid/reduce escaping delimiter character
<code>\t</code>	escape sequences like those supported in string literals
	but escapes like <code>\b</code> and <code>\s</code> have different meaning in regexps
<code>\%</code>	undefined escapes will match the character it escapes
<code>\x7c</code>	will match <code> </code> literally
	instead of acting as alternation metacharacter

Exercises

a) Transform given input strings to expected output using same logic on both strings.

```
>> str1 = '(9-2)*5+qty/3'
>> str2 = '(qty+4)/2-(9-2)*5+pq/4'

>> str1.gsub()      ##### add your solution here
=> "35+qty/3"
>> str2.gsub()      ##### add your solution here
=> "(qty+4)/2-35+pq/4"
```

b) Replace `(4)\|` with `2` only at the start or end of given input strings.

```
>> s1 = '2.3/(4)\|6 foo 5.3-(4)\|'
>> s2 = '(4)\|42 - (4)\|3'
>> s3 = "two - (4)\|\n"

>> pat =      ##### add your solution here

>> s1.gsub(pat, '2')
=> "2.3/(4)\|6 foo 5.3-2"
>> s2.gsub(pat, '2')
=> "242 - (4)\|3"
>> s3.gsub(pat, '2')
=> "two - (4)\|\n"
```

c) Replace any matching item from given array with `X` for given input strings. Match the elements from `items` literally. Assume no two elements of `items` will result in any matching conflict.

```
>> items = ['a.b', '3+n', 'x\y\z', 'qty||price', '{n}']

>> pat =      ##### add your solution here

>> '0a.bcd'.gsub(pat, 'X')
=> "0Xcd"
>> 'E{n}AMPLE'.gsub(pat, 'X')
=> "EXAMPLE"
>> '43+n2 ax\y\ze'.gsub(pat, 'X')
=> "4X2 aXe"
```

d) Replace backspace character `\b` with a single space character for the given input string.

```
>> ip = "123\b456"
>> puts ip
12456

>> ip.gsub()      ##### add your solution here
=> "123 456"
```

e) Replace all occurrences of `\o` with `o`.

```
>> ip = 'there are c\omm\on aspects am\ong the alternati\ons'

>> ip.gsub()      ##### add your solution here
=> "there are common aspects among the alternations"
```

f) Replace any matching item from the array `eqns` with `X` for given the string `ip` . Match the items from `eqns` literally.

```
>> ip = '3-(a^b)+2*(a^b)-(a/b)+3'
>> eqns = %w[(a^b) (a/b) (a^b)+2]

>> pat =          ##### add your solution here

>> ip.gsub(pat, 'X')
=> "3-X*X-X+3"
```