

# Assignment Project Exam Help

Introduction to Regular Expressions

<https://powcoder.com>

Faculty of Information Technology  
Monash University

Add WeChat FIT5196 week 2 powcoder

## Regular Expressions

- A regular expression is a set of symbols that describes a text pattern.

- ▶ `\d{4}-\d{2}-\d{2}`
- ▶ `wrangling`

# Assignment Project Exam Help

- Why regular expressions?

- ▶ Regular expressions are useful in finding, replacing and extracting information from text, such as log files, HTML/XML files, and other documents

- ▶ Search a document for color or neighborhood with or without a

- Convert a tab-delimited file to a comma-delimited file
- Find duplicated words in a text
- Search and replace “Bob” and “Bobby” with “Robert”

- ▶ Regular expressions are useful in verifying whether input fits into the text pattern, such as verifying

- phone numbers: Does a phone number have the right number of digits?
- emails: Is an email address in a valid format?
- date: Is a date in the right format? Does the month exceed 12?

## Regular Expressions: validate emails

`r"([a-zA-Z0-9_+-.]+@[a-zA-Z0-9]+\.[a-zA-Z0-9-]+\.)+$)"1`

# Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

---

<sup>1</sup><http://emailregex.com/>

## Regular Expressions: validate emails

`r"([a-zA-Z0-9_+-.]+@[a-zA-Z0-9]+\.[a-zA-Z0-9-]+)$"`<sup>1</sup>

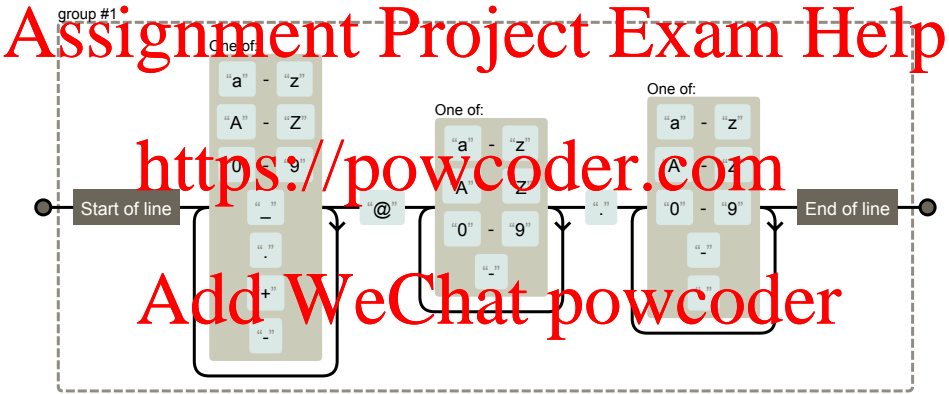


Figure: Figure generated by <https://regexper.com/>

<sup>1</sup><http://emailregex.com/>

### 1 Regular Expression Syntax

- character sets
- repetition
- grouping
- raw string in Python

<https://powcoder.com>

### 2 Cases studies

Add WeChat powcoder

### 3 Summary

## Matching String Literals

The most obvious feature of regular expressions is matching strings with one or more literal characters, called string literals.

- Everything is essentially a character in regular expressions
  - ▶ `cat` matches "cat".
  - ▶ `cat` matches the first three characters of "cattle" and "catfish".
- It is similar to searching in word processing program
- Matching is case sensitive.
  - ▶ `cat` does not match "Cat".
- How does regular expression engine work?

**Add WeChat powcoder**

The cow, camel and cat communicated.

## Character sets: [...]

Assume that we are going to match the following two words:

Assignment <sup>grey gray</sup> Project Exam Help

What the regular expression should be?

<https://powcoder.com>

Add WeChat powcoder



## Character sets: [...]

Assume that we are going to match the following two words:

Assignment Project Exam Help

grey gray

What the regular expression should be?

- [...] indicate a set of characters
  - ▶ Matches any one of several characters in the set, but only one
  - ▶ The order of characters does not matter.

<https://powcoder.com>

Add WeChat powcoder



## Character sets: [...]

Assume that we are going to match the following two words:

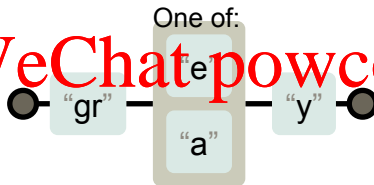
# Assignment Project Exam Help

What the regular expression should be?

- [...] indicate a set of characters
  - ▶ Matches any one of several characters in the set, but only one
  - ▶ The order of characters does not matter.
  - ▶ The regular expression is `gr[ea]y`

One of:

# Add WeChat powcoder



- ▶ `gr[ea]y` does not match `grAy`, `graay`, and `graey`.



**Character ranges:**  $[a - zA - Z]$  and  $[0 - 9]$

Assume that we are going to match victory car plate numbers, for example

# Assignment Project Exam Help

XRA-000 1AA 1AA

Note the letters can be from A to Z, and the numbers can be from 0 to 9. What the regular expression should be?

<https://powcoder.com>

Add WeChat powcoder

**Character ranges:**  $[a - zA - Z]$  and  $[0 - 9]$

Assume that we are going to match victory car plate numbers, for example

# Assignment Project Exam Help

XRA 000 1AA 1AA

Note the letters can be from A to Z, and the numbers can be from 0 to 9. What the regular expression should be?

- Character ranges can be indicated by giving two characters and separating them by a '-'.  
<https://powcoder.com>

- Example:

- $[0 - 9]$
- $[a - z]$  or  $[A - Z]$

- Caution

- $[50 - 99]$  is not all numbers from 50 to 99, it is the same as  $[0 - 9]$ .

One of:



**Character ranges:**  $[a - zA - Z]$  and  $[0 - 9]$

Assume that we are going to match victory car plate numbers, for example

XRA 000 1AA 1AA

# Assignment Project Exam Help

Note the letters can be from A to Z, and the numbers can be from 0 to 9. What the regular expression should be?

- $[A-Z0-9][A-Z][A-Z]\backslashs[0-9][A-Z0-9][A-Z0-9]$

<https://powcoder.com>





## Negative character sets: `[^...]`

Assume that we are going to write a regular expression that matches only the live animals

# Assignment Project Exam Help

Question: what is the regular expression?

<https://powcoder.com>

Add WeChat powcoder



## Negative character sets: `[^ ...]`

Assume that we are going to write a regular expression that matches only the live animals

# Assignment Project Exam Help

Question: what is the regular expression?

- `[^ ...]`: If the first character of the set is `^`, all the characters that are not in the set will be matched.

- ▶ `[^h]og` matches "hog" and "dog", but not "log".

- ▶ Caution:

- Does `see[^mn]` match "see"?

- Does `see[^mn]` match "see "?

<https://powcoder.com>  
Try the regular expression in Pythex (<http://pythex.org/>)!



## Metacharacters inside character sets: `[.+]`

Assume that we are going to match the following two strings:

# Assignment Project Exam Help

Now, we need to match `()` and `[]`. How can we do that?

<https://powcoder.com>

Add WeChat powcoder



## Metacharacters inside character sets: `[.+]`

Assume that we are going to match the following two strings:

# Assignment Project Exam Help

Now, we need to match `()` and `[]`. How can we do that?

- Metacharacters inside character sets are already escaped. In other words they lose their special meaning inside sets

- ▶ Example

- `h[ai.u]t` matches “hat”, “h.t”, but not “hot”

- Exceptions: `]`, `-`, `^` and `\` that do need to be escaped.

- ▶ `h[a.u]t` → `h[ai.]t`?

# Add WeChat powcoder



## Metacharacters inside character sets: `[.+]`

Assume that we are going to match the following two strings:

# Assignment Project Exam Help

Now, we need to match `()` and `[]`. How can we do that?

- Metacharacters inside character sets are already escaped. In other words they lose their special meaning inside sets.

▶ Example

– `h[ai.ut]` matches `hat`, `h.t`, but not `hot`

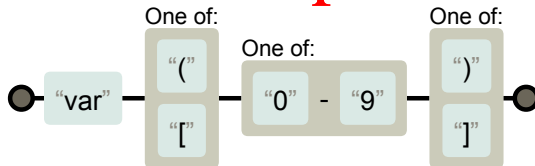
- Exceptions: `]`, `-`, `^` and `\` that do need to be escaped.

▶ `h[ai.ut]` → `h[ai\ut]`?

- `var([(){}9])`

<https://powcoder.com>

Add WeChat powcoder



## Metacharacters inside character sets: `[.+]`

Assume that we are going to match the following two strings:

# Assignment Project Exam Help

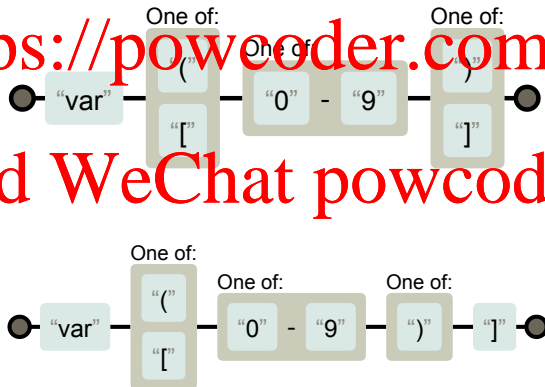
Now, we need to match `()` and `[]`. How can we do that?

- `var([0-9]\d\)`

# https://powcoder.com

# Add WeChat powcoder

- `var([0-9]\d\)`



## Shorthand character sets

Shorthand	meaning	Equivalent
<code>\d</code>	matches any decimal digit from 0 to 9	<code>[0-9]</code>
<code>\w</code>	matches any word character	<code>[a-zA-Z0-9_]</code>
<code>\s</code>	matches any white space character	<code>[\t\n\r]</code>
<code>\D</code>	matches any non-digit character;	<code>[^0-9]</code>
<code>\W</code>	matches any non-alphanumeric character	<code>[^a-zA-Z0-9_]</code>
<code>\S</code>	matches any non-whitespace character	<code>[^\t\n\r]</code>

Add WeChat powcoder



## Shorthand character sets

Shorthand	meaning	Equivalent
<code>\d</code>	matches any decimal digit from 0 to 9	<code>[0-9]</code>
<code>\w</code>	matches any word character	<code>[a-zA-Z0-9_]</code>
<code>\s</code>	matches any white space character	<code>[\t\n\r]</code>
<code>\D</code>	matches any non-digit character;	<code>[^0-9]</code>
<code>\W</code>	matches any non-alphanumeric character	<code>[^a-zA-Z0-9_]</code>
<code>\S</code>	matches any non-whitespace character	<code>[^\t\n\r]</code>

### Examples:

- ▶ `\d\d\d\d` matches four digit numbers, such as "2018", but not text.
- ▶ `\w\w\w` matches three word characters such as "abc", "123" and "d\_b"
- ▶ `\w\w\s\w` matches "ab c" but not "a bc".
- ▶ `[\w]-[\w]` matches two characters separated by a hyphen.
- ▶ `[^\d]` is the same as `[\D]`

## Shorthand character sets

- Caution:

- ▶ Is `[\^d\s]` the same as `[\D\S]`?

# Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

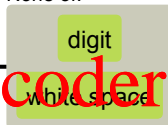
## Shorthand character sets

- Caution:

- Is `[^\d\s]` the same as `[\D\S]`?

- `[^\d\s]` Not digit OR space character

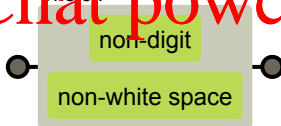
None of:



<https://powcoder.com>

- `[\D\S]` EITHER NOT digit OR NOT space character

One of:



Add WeChat powcoder

Try the regular expression with the following sentence: "Data Wrangling S2  
2018 week 2"

## Repetition Expressions: repetition meta-characters

meta-characters	meaning
*	Match 0 or more repetitions of the preceding regex
+	Match 1 or more repetitions of the preceding regex
?	Match 0 or 1 repetitions of the preceding regex

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## Repetition Expressions: repetition meta-characters

meta-characters	meaning
*	Match 0 or more repetitions of the preceding regex
+	Match 1 or more repetitions of the preceding regex
?	Match 0 or 1 repetitions of the preceding regex

### Examples:

- Assume we are going to match the following words

<https://powcoder.com>

oops oops oooooops oooooops

but not

Add WeChat oops powcoder

which regular expression(s) should we use?

- 1 oo\*ps
- 2 ooo\*ps
- 3 oo+ps
- 4 oo?ps



## Repetition Expressions: repetition meta-characters

meta-characters	meaning
*	Match 0 or more repetitions of the preceding regex
+	Match 1 or more repetitions of the preceding regex
?	Match 0 or 1 repetitions of the preceding regex

### Examples:

- Assume we are going to match the following words

<https://powcoder.com>

oops oops oooooops oooooops

but not

Add WeChat oops powcoder

which regular expression(s) should we use?

- 1 oo\*ps
- 2 ooo\*ps
- 3 oo+ps
- 4 oo?ps



## Repetition Expressions: repetition meta-characters

meta-characters	meaning
*	Match 0 or more repetitions of the preceding regex
+	Match 1 or more repetitions of the preceding regex
?	Match 0 or 1 repetitions of the preceding regex

### Examples:

- Assume we are going to match the following words

<https://powcoder.com>

oops oops oooooops oooooops

but not

Add WeChat <sup>oops</sup> powcoder

which regular expression(s) should we use?

- 1 oo\*ps
- 2 ooo\*ps
- 3 oo+ps
- 4 oo?ps



## Repetition Expressions: repetition meta-characters

meta-characters	meaning
*	Match 0 or more repetitions of the preceding regex
+	Match 1 or more repetitions of the preceding regex
?	Match 0 or 1 repetitions of the preceding regex

### Examples:

- Assume we are going to match the following words

<https://powcoder.com>  
oops oops oooooops oooooops

but not

Add WeChat oops  
which regular expression(s) should we use?

- 1 oo\*ps
- 2 ooo\*ps
- 3 oo+ps
- 4 oo?ps



## Repetition Expressions: repetition meta-characters

meta-characters	meaning
*	Match 0 or more repetitions of the preceding regex
+	Match 1 or more repetitions of the preceding regex
?	Match 0 or 1 repetitions of the preceding regex

### Examples:

- Assume we are going to match the following words

<https://powcoder.com>  
oops oops oooooops oooooops

but not

oops  
Add WeChat powcoder  
which regular expression(s) should we use?

- 1 oo\*ps
- 2 ooo\*ps
- 3 oo+ps
- 4 oo?ps



## Repetition Expressions: repetition meta-characters

meta-characters	meaning
*	Match 0 or more repetitions of the preceding regex
+	Match 1 or more repetitions of the preceding regex
?	Match 0 or 1 repetitions of the preceding regex

- Examples:

- Assume we are going to match the following words

<https://powcoder.com>

oops oops oooooops ooooooops

but not

Add WeChat oops powcoder

The regular expressions that we can use:

`ooo*ps`   `oo+ps`

Try the regular expression in Pythex!

## Repetition Expressions: quantified repetitions

- `{m, n}`: matches exactly from  $m$  to  $n$  repetitions of the preceding regular expression.

- ▶  $m$  (min) and  $n$  (max) are positive numbers
- ▶  $m$  must be always be included, can be 0
- ▶  $n$  is optional

- Three syntax

- ▶ `\d{2}` matches numbers with exactly 2 digits.
- ▶ `\d{2, 4}` matches numbers with 2 to 4 digits.
- ▶ `\d{2, }` matches numbers with at least 2 digits ( $n$  is infinite).

Try the 'loops' example in Pythex, but with `{m, n}`

## Repetition Expressions: quantified repetitions

- Suppose we are going to match the following

report_2018_09		assignment_2018_9
budget_18_08		assignment_18_7

but not

report_201809_39		assignment_8_9000
budget_2345678_08		assignment_000999_7

what is the regular expression?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## Repetition Expressions: quantified repetitions

- Suppose we are going to match the following

report_2018_09		assignment_2018_9
budget_8_08		assignment_18_7

but not

report_201809_39		assignment_8_9000
budget_2345678_08		assignment_000999_7

what is the regular expression?

`\w+_d{2,4}_d{1,2}`



Try the regular expression in Pythex



## Repetition Expressions: greedy v.s. lazy regex

- Greedy strategy:

- ▶ Match as much as possible before giving control to the next regular expression part.  
Regular expressions try to match the longest possible string

- ▶ Examples

`.*\d+`  
number 816  
<https://powcoder.com>

Add WeChat powcoder

## Repetition Expressions: greedy v.s. lazy regex

- Greedy strategy:

- Match as much as possible before giving control to the next regular expression part.
- Regular expressions try to match the longest possible string

- Examples

- Question: Given a string like `.*\d+`  
`https://powcoder.com number 816`

"data", "wrangling", "FIT5196, S2."

What is the match of regular expression `"[a-z]+"`?

- "data", "wrangling"
- "wrangling", "FIT5196, S2."
- "data", "wrangling", "FIT5196, S2."

## Repetition Expressions: greedy v.s. lazy regex

- Greedy strategy:

- ▶ Match as much as possible before giving control to the next regular expression part.  
Regular expressions try to match the longest possible string

- ▶ Examples

- ▶ `.*\d+`  
number 816  
<https://powcoder.com>

- ▶ Question: Given a string like

"data", "wrangling", "FIT5196, S2."

what is the match of regular expression `".*", "\d+"`?

- 1 "data", "wrangling"
- 2 "wrangling", "FIT5196, S2."
- 3 "data", "wrangling", "FIT5196, S2."

## Repetition Expressions: greedy v.s. lazy regex

- Lazy strategy:

- ▶ Match as little as possible before giving control to the next regular expression part

- ▶ Syntax

- `*?`

- `+?`

- `??`

- `{m,n}?`

- ▶ Example:

`.*?\d+`

number 516

Add WeChat powcoder

## Repetition Expressions: greedy v.s. lazy regex

- Lazy strategy:

- ▶ Match as little as possible before giving control to the next regular expression

- ▶ part  
Syntax

- `*?`

- `+?`

- `??`

- `{n,m}?`

- ▶ Example:

`.*?\d+`

number 516

- Question: Given a string like

"data", "wrangling", "FIT5196, S2."

what is the match of regular expression `".+?"`, `".+?"`?

- 1 "data", "wrangling"
- 2 "wrangling", "FIT5196, S2."
- 3 "data", "wrangling", "FIT5196, S2."

## Repetition Expressions: greedy v.s. lazy regex

- Lazy strategy:

- ▶ Match as little as possible before giving control to the next regular expression

- ▶ part  
Syntax

- `*?`

- `+?`

- `??`

- `{n,m}?`

- ▶ Example:

`.*?\d+`

number 516

- Question: Given a string like

"data", "wrangling", "FIT5196, S2."

what is the match of regular expression `".+?"`, `".+?"`?

- 1 "data", "wrangling"
- 2 "wrangling", "FIT5196, S2."
- 3 "data", "wrangling", "FIT5196, S2."

## Grouping: ( ... )

- ( ... ) matches whatever regular expression is inside the parentheses, and indicates the start and end of a group.

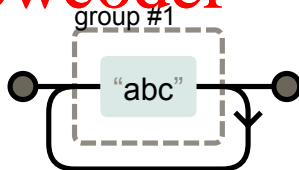
- ▶ Apply repetition operators to a group of regular expressions
- ▶ Makes regular expressions easier to read
- ▶ Capture groups for use in matching, replacing and extraction, i.e., the contents of a group can be retrieved.
- ▶ Cannot be used inside a character set

- examples:

- ▶ `abc` matches `abc`, `abcc`, `abccc`



- ▶ `(abc)+` matches `abc`, `abcbcb`, `abcbcbcbcb`



## Grouping: (...)

- (...) matches whatever regular expression is inside the parentheses, and indicates the start and end of a group.

- ▶ Apply repetition operators to a group of regular expressions
- ▶ Makes regular expressions easier to read
- ▶ Capture groups for use in matching, replacing and extraction, i.e., the contents of a group can be retrieved.
- ▶ Cannot be used inside a character set

- examples:

- ▶ "Incident American Airlines Flight 11 involving a Boeing 767-223ER in 2001"
- ▶ Regular expression: `Incident (.*) involving`



- ▶ Try it with python script!!!





## Alternation: |

- “|” is an OR operator

- ▶  $A|B$  will match any string that matches either A or B
- ▶ Ordered – left most expression gets precedence.
- ▶ Multiple patterns can be daisy-chained.
- ▶ Group alternation expressions to keep them distinct.

- Examples:

- ▶  $apple|orange$  matches “apple” and “orange”
- ▶  $(apple|orange) juice$  matches “apple juice” and “orange juice”
- ▶  $w(ei|ie)rd$  matches both “weird” and “wierd”.

# Assignment Project Exam Help

<https://powcoder.com>

## Add WeChat powcoder

## The backslash plague: \

- The back slash \ indicates special forms or to allow special characters to be used without invoking their special meaning.<sup>2</sup>

Characters	Stage
<code>\section</code>	text string to be matched
<code>\\section</code>	Escaped backslash for <code>re.compile()</code>
<code>\\\\section</code>	Escaped backslashes for a Python string literal

- So, to match a literal backslash, one has to write `\\\\` as the regular expression string
- Can we simply the expression?

<sup>2</sup>see <https://docs.python.org/3/howto/regex.html>

## Raw String: `r"... "`

- Raw String suppress actual meaning of escape characters, and do not treat the backslash as a special character at all.

Assignment Project Exam Help

Regular Python string literal	Raw string
<code>"\\\\section"</code>	<code>r"\\section"</code>
<code>"\\w+\\s+"</code>	<code>r"\\w+\\s+"</code>

- Regular expressions will often be written in Python code using this raw string notation.
- Try the Python script!!!

<https://powcoder.com>

Add WeChat powcoder

## Case study 1: validate dates

- Date samples (day, month, year):

- ▶ 02/08/2018
- ▶ 2/8/2018
- ▶ 2/8/18
- ▶ 23/08/2018
- ▶ 23-08-2018

- See <https://www.powcoder.com>

Add WeChat powcoder

## Case study 2: validate credit card number

- Assume that you're given the job of implementing an order form for a company that accepts payment by credit card issued by the world's major credit card companies such as VISA, MasterCard, and American Express.

- ▶ All Visa card numbers start with a 4. New cards have 16 digits. Old cards have 13.
  - 4123456789012
  - 4123456789012345
- ▶ MasterCard numbers either start with the numbers 51 through 55. All have 16 digits.
  - 5123456789012345
  - 5523456700012345
- ▶ American Express card numbers start with 34 or 37 and have 15 digits.
  - 341234567890123
  - 371234567890123

- See jupyter notebook!!!

## Summary: what to do this week

- Regular expressions are the major tool used in data parsing.
- Study materials provided in Moodle.
- Attend tutorial 2.
  - ▶ Try to finish all the materials provided in the tutorial.
- **Assessment 1** will be released at the end of this week (week 2).
- Topic for next week:
  - ▶ Parsing data stored in different file formats, CSV, JSON, XML, EXCEL, and PDF

# Assignment Project Exam Help

<https://powcoder.com>

## Add WeChat powcoder