# Regex For Dummies. Part 4: Capturing Groups and Backreferences

NALSengineering · Follow

7 min read · Sep 27, 2023

Listen    Share    More

See <u>Part 1 here</u>: Quantifiers

See <u>Part 2 here</u>: Flavors, Flags, and Assertions

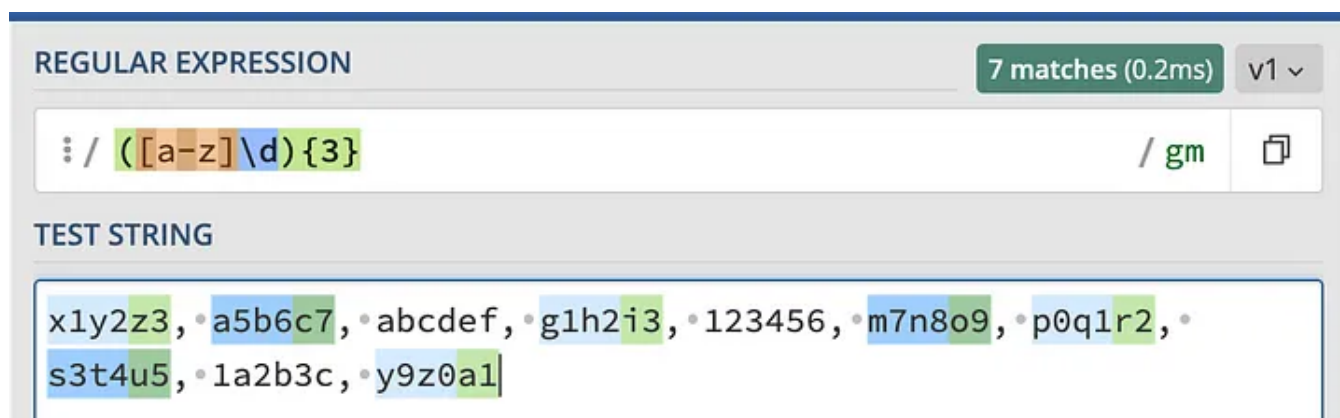See <u>Part 3 here</u>: Character Set, Or condition, and Word Boundary Assertion

In this part, we will find out the advanced concept in Regex: Capturing Groups and Backreferences

### 1. **Capturing Groups** `()`

Capturing groups allow you to treat a part of your regex pattern as a single unit. This is especially useful when you want to apply quantifiers or modifiers to multiple characters or subpatterns. For example, `(abc)+` matches one or more repetitions of "abc" as a whole.

Suppose you have a set of IDs like `x1y2z3`, `a5b6c7`, `abcdef`, `g1h2i3`, `123456`, `m7n8o9`, `p0q1r2`, `s3t4u5`, `1a2b3c`, `y9z0a1`, and you want to find out which of them follow the pattern `letter number letter number letter number`. There are two ways to do that:

- Without using grouping: `[a-z]\d[a-z]\d[a-z]\d`

- Using grouping, you can make the pattern shorter by grouping the `[a-z]\d` sequence and applying an exact quantifier of `3` to it: `([a-z]\d){3}`



Capturing groups are very useful for extracting data.

Let's say we want to extract the day, month, and year from birthday dates in the format `dd-mm-yyyy`, we can use this regex:

```
(\d{1,2})-(\d{1,2})-(\d{4})
```

We have three capturing groups: the 1st group captures the day, the 2nd group captures the month, and the 3rd group captures the year (LINK).

I'll teach you how to extract them by coding in the next part.

## 2. Backreferences

Backreferences are a feature that allows you to reference and reuse the text captured by a capturing group within the same regex pattern. They provide a powerful way to search for repeated patterns and validate complex text structures.

In a regex pattern, you can reference a capturing group by using a backslash followed by the group number. Group numbers are assigned based on the order of opening parentheses in the regex pattern, starting from 1. For example, `\1` refers to the text captured by the first capturing group, `\2` to the second, and so on.

Let's explore backreferences with an example:

- `(abc)\1` matches `abcabc`. In which, `(abc)` is a capturing group, and `\1` is a backreference that matches the same text as captured by the capturing group, so, `\1` matches `abc` too.

- `(a|b|cd)\1` matches `aa`, `bb` or `cdcd`

- `(\d)(.+)\1+\2?` matches `1a1`, `1a1a`, `1a11`, `1a11a1`

You can also easily check the captured texts of the groups on the regex101.



Capturing groups and backreferences are very useful in many scenarios, such as the "finding of duplicated words" problem.

Suppose we need to search for duplicated words in this article, such as "the the regex" or "hello hello world". Imagine how challenging this task would be if you were to use the regular search. However, we can easily accomplish it with the regex `\b([A-Za-z]+) +\1\b` (LINK).



In which:

- `\b` : assert the start position of a word.

- `([A-Za-z]+)` : This is the first capturing group. It captures one or more alphabetical characters (both uppercase and lowercase).

- `+` : This part matches one or more spaces. It allows for the possibility of spaces between the duplicated words.

- `\1` : This is a backreference to the first capturing group. It checks that the same word captured by the first group appears again immediately after the spaces.

- `\b` : assert the end position of a word.

## 3. Non-capturing groups

It is also possible to make a group non-capturing. That way, you won't be able to refer to it in the pattern. To create a **non-capturing** group, you use a question mark and a colon right after the opening parenthesis.

The syntax for that looks like this:

```
(?:chars)
```

For instance, in the example "extracting the day, month, and year from birthday dates" above, if we don't want to capture the first two groups, then we would modify the regex as follows:

```
(?:\d{1,2}-){2}(\d{4}):\1
```

I added the backreference `\1` to test the 1st capturing group.

## 4. Nested capturing groups

For nested capturing groups, the order of the groups is the same as the order of the left parentheses. For example, in this regex `(\d{4})-((\d{3})-(\d{2}))`,

- `(\d{4})` is the first capturing group and is numbered 1.

- `((\d{3})-(\d{2}))` is the second capturing group and is numbered 2.

- `(\d{3})` is the first inner capturing group of the second group and is numbered 3.

- `(\d{2})` is the second inner capturing group of the second group and is numbered 4.



## 5. Named capturing group

Named capturing groups allow you to assign names to your capturing groups, making it easier to reference and work with specific matched portions of text. Instead of referring to capturing groups by their numerical indices, you can use descriptive names, which enhances the readability and maintainability of your regex patterns.

To define a named capturing group, you use the `(?<name> ... )` syntax, where `<name>` is the name you want to assign to the group. Here's an example:

```
(?<year>\d{4})-(?<month>\d{2})-(?<day>\d{2})
```

In this pattern:

- `(?<year>\d{4})` defines a named capturing group named "year" to capture a 4-digit year.

- `(?<month>\d{2})` defines a named capturing group named "month" to capture a 2-digit month.

- `(?<day>\d{2})` defines a named capturing group named "day" to capture a 2-digit day.



Named capturing groups

## 6. Find and Replace with Regex in Text Editors

Most Text Editors such as Google Docs, Google Sheets, Google Slides, and IDEs like VSCode and Android Studio support the Find and Replace with Regex feature.

### 6.1. Google Sheets

Let's say we need to find and replace dates in the format `dd-mm-yyyy` with `dd/mm/yyyy`. We can easily do this using the regex:

```
(\d{1,2})-(\d{1,2})-(\d{4})
```

| | A |
|---|---|
| 1 | 31-08-1994 |
| 2 | 1-6-2020 |
| 3 | 5-12-1997 |
| 4 | 13-11-1998 |
| 5 | |
| 6 | |
| 7 | |
| 8 | |
| 9 | |
| 10 | |
| 11 | |
| 12 | |
| 13 | |
| 14 | |
| 15 | |
| 16 | |
| 17 | |
| 18 | |
| 19 | |
| 20 | |
| 21 | |
| 22 | |
| 23 | |

**Find and replace**                                                    ✕

Find            $(\d{1,2})-(\d{1,2})-(\d{4})$

Replace with    $1/$2/$3

Search          All sheets ▾

☐  Match case

☐  Match entire cell contents

☑  Search using regular expressions   Help

☐  Also search within formulas

☐  Also search within links

Find      Replace      Replace all      Done

Please note that we need to enable the "Search using regular expressions" mode.

Search          All sheets ▾

☐  Match case

☐  Match entire cell contents

☑  Search using regular expressions   Help

☐  Also search within formulas

☐  Also search within links

## 6.2. Google Docs

Let's say we want to search for duplicated words in the Google Docs document. We could use the regex `\b([A-Za-z]+) +\1\b`. However, according to the official documentation, capture groups only work with Google Sheets. That means we cannot use replacement backreferences in Google Docs but have to do it manually.

It was a a beautiful day, and the sun was shining shining brightly. I decided to to take a walk in the park to to enjoy the weather. As I walked, I saw a a group of children playing playing in the playground. They were having so much fun, running running around and laughing laughing loudly. I also saw a a couple sitting sitting on a bench, holding hands and gazing gazing at the lake. It was a a peaceful scene, and I couldn't help but smile as I continued my walk walk.|

**Find and replace**                                              ✕

Find          `\b([A-Za-z]+) +\1\b|`                      1 of 13

Replace with  [                                              ]

☐ Match case

☑ Use regular expressions (e.g. \n for newline, \t for tab) Help

☑ Ignore Latin diacritics (e.g. ä = a, E = É)

[ Replace ]   [ Replace all ]   [ Previous ]   [ Next ]

## 6.3. VSCode IDE

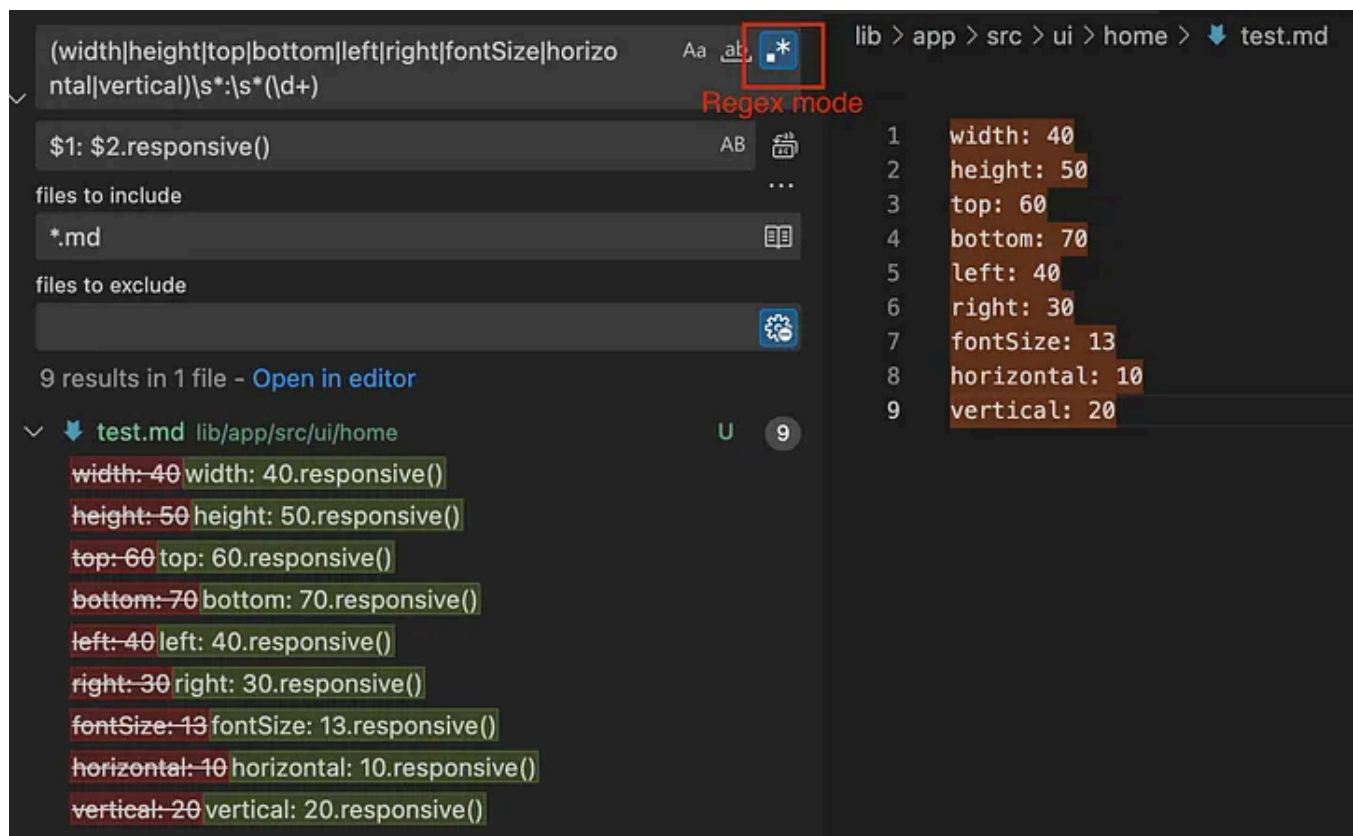Let's get back to the problem introduced in the first part.

"Have you ever been involved in a project with over 500 code files and realized there are a lot of hard-coded dimensions? Your task is to find and fix them like this"

```
width: 40        ->   width: 40.responsive()
height: 50       ->   height: 50.responsive()
top: 60          ->   top: 60.responsive()
bottom: 70       ->   bottom: 70.responsive()
```

```
left: 40        ->  left: 40.responsive()
right: 30       ->  right: 30.responsive()
fontSize: 13    ->  fontSize: 13.responsive()
horizontal: 10  ->  horizontal: 10.responsive()
vertical: 20    ->  vertical: 20.responsive()
```

This is a big problem if we don't know how to use Regex. But now we can deal with it with just one click.

Don't forget to enable the "Regex mode".



Looking into the regex:

```
(width|height|top|bottom|left|right|fontSize|horizontal|vertical)\s*:\s*(\d+)
```

- `(width|height|top|bottom|left|right|fontSize|horizontal|vertical)` : This is the first capturing group. It captures the property name, such as width, height,...

- `\s*` represents zero or more whitespace characters.

- `(\d+)` : This is the second capturing group. It captures the dimension value.

- `$1` : This is the backreference of the first capturing group.

- `$2` : This is the backreference of the second capturing group.

## Conclusion

We have been able to apply regex to our work so far. In the next part, I will introduce you to an advanced feature in regex: Lookaround Groups, which include Lookaheads and Lookbehinds.

Continue to Part 5:

**Regex for Dummies. Part 5: Lookaround Assertions — Lookaheads and Lookbehinds**

Credit: Nguyễn Thành Minh (Android Developer)

medium.com

Regex · Android · IOS · Flutter · Java

Follow

## Written by NALSengineering

1.3K Followers

Knowledge sharing empowers digital transformation and self-development in the daily practices of software development at NAL Solutions.