

How to use 'Regular Expressions' to speed up the corpus building process

A guide for the Oslo Medical Corpus

Purpose: A 'regular expression', or 'regex' for short, is a sequence of special text characters which can be used to describe a specific re-occurring pattern within a body of text.

By using regular expressions, the process of preparing and annotating texts for uploading to a corpus can be made much less time- and labour-intensive.

This document provides a brief introduction to regular expressions and the ways in which they can be helpful for the purposes of corpus-building. Additional information can be found in online tutorials such as <http://www.regular-expressions.info/>.

Regular expressions can be applied in text editing software such as *Visual Studio Code*. They are used to enhance the Search and Replace function (Ctrl + F) of this software and can help to identify and manipulate certain frequently occurring features of a digitised text.

Literal and special characters

When using regular expressions, there are two kinds of character: 'literal' characters and 'special' characters.

All alphanumeric characters (the letters **a-z** and the numbers **0-9**) and some punctuation marks are 'literal' characters. Literal characters do not normally have any special function within the regular expression.

For example, if you type the word **evidence** into the search box, the software will simply find every instance in the text where these nine literal characters (**e, v, i, d, e, n, c, e**) are found next to each other (and in this order) within the body of the text.

'Special' characters, on the other hand, have a particular function within the regular expressions syntax.

Below are the special characters and their functions:

1. period or dot .

The dot matches any (single) character, including blank spaces but except line break characters. It can be used in much the same way as a 'wildcard' token used in other search engine software.

For example, typing **w.nt** into the search box matches **went, want, wont**, etc., as well as character sequences such as **w-nt**.

This can be useful for checking consistency in spelling, and for identifying and fixing OCR errors in case the text has been scanned and converted into a text file.

2. the question mark ?

Placing a question mark after a character renders it optional.

For example, searching for **honou?r**, matches both **honour** and **honor**.

3. square brackets []

Square brackets can be placed in a search string to designate a 'character class'.

For example, **[a-z]** will match any lower-case letter of the alphabet, **[A-Z]** will match any upper-case letter of the alphabet and **[a-zA-Z]** will match any letter of the alphabet, including both upper- and lower-case characters.

[0-9] matches any numerical digit.

4. curly braces {}

If you want to search for and find more than one character in a character class, then you must specify this using curly braces.

For example, **[0-9]{2}** will match any sequence of two digits, i.e. **12** or **98** or **45**.

You can specify a minimum and a maximum using a comma to separate: **{min,max}**.

For example, **[0-9]{2,}** will match any sequence of at least two digits, while **[0-9]{1,3}** will match any sequence of numbers between one and three digits long, i.e. **5** or **98** or **404**. This can be useful for finding and removing page numbers placed at the top/bottom of each page of a document.

5. the plus sign +

The plus sign matches the preceding character or character class when it appears one or more times.

For example, searching for **a+** will find **a**, **aaaa**, **aaaaaaaaaaaaa** and **aaaaaaaaaaaaaaaaa**.

[a-z]+ will find **a**, **aaaaaaaaaaaaa**, **bbbbbbbb** and **greaagdgadfdharhgbvxcvbbttttb**

6. the asterisk or star *

The asterisk is used to find instances where a specific character or character class is repeated any number of times, including zero times.

For example, searching for **conv*** will find **conversing**, but also **consider**.

7. the caret ^

If you want to omit ('negate') a character or character class from your search, you can place the ^ symbol before this character or character class in the square brackets.

For example, **q[^u]** matches any **q** and the character that follows it (including a white space), as long as this second character is not a **u**.

[^a-z] matches any character (including a white space) which is not a lower-case letter of the alphabet.

8. the vertical bar or pipe symbol |

The vertical bar means 'EITHER/OR'.

For example, if you want a regex that matches either **online discussion** or **on-line discussion**, you should search for **(online|on-line) discussion**. The role of the parentheses is outlined below.

9. the backslash \

The backslash can be used if you want to cancel or 'escape' any special character's special function.

For example, if you want to search the text for a full stop, you must place a backslash in front of this character **\.**, otherwise the software will treat the full stop as a special character and match every single character in the document.

The backslash is also used to give certain literal characters 'special' functions.

\n matches a new line or 'line break';

\t matches a tabbed space.

\s matches any white space, including tabs and line breaks;

\w matches a 'word character' (alphanumeric characters plus underscore);

\d matches a single digit.

10. parentheses () and the dollar sign \$

Parentheses are used to 'group' certain characters.

For instance, in the example given above for the vertical line, we used parentheses to tell the search engine to 'group' **online** and **on-line** as alternatives in the EITHER/OR query.

In combination with the dollar sign, they can also be used to 'grab' a certain set of characters and re-use these same characters via the replace function.

The dollar sign is used to represent any text that you have ‘grabbed’ in the search string and that you want to reuse in the replace string. Grabbed sections are numbered left to right in sequence: **\$1**, **\$2**, **\$3**, **\$4**...

See the example given below regarding hyphens for an illustration of how to use parentheses and dollar signs.

Applying regular expressions to corpus text preparation

Regular expressions cannot be used to solve every kind of problem in the process of text annotation. However, they are particularly useful for certain kinds of operation. Below are a series of concrete examples to illustrate the ways in which regular expressions can be applied to corpus text preparation. Please note that every document is different so these regular expressions will likely need to be adapted to each particular case.

- *Paragraphing:*

One of the most common and simplest uses for regular expressions is for removing ‘white space’ and blank lines from the xml document. For example, if after every paragraph in the document there are two (or more) line breaks instead of just one, you can simply Search **\n{2,}** and Replace this with **\n**.

Likewise, tabbed spaces at the beginning of each paragraph can be quickly removed by searching **\n\t** and replacing this with **\n**.

- *Removing hyphens in the middle of a word:*

Some publishers may regularly split and hyphenate words located at the end of a line in order to improve the layout of the text on the printed page. For the corpus file, however, these hyphens will need to be removed and regular expressions provide an easy way of doing this. Simply search for **(([a-z])- ([a-z]))** and replace this with **\$1\$2**. This will remove both the hyphen and the space after the hyphen from the middle of the word: **demo- cracy** becomes **democracy**.

- *Identifying web addresses:*

Since all urls follow a very similar pattern, these can easily be found using regular expressions. Search for **https?://[^\<>]+** and, in case such instances need to be removed, replace with nothing.

- *Headers and page numbers:*

In some documents, particularly scanned copies of physical books, you will potentially find a header (and perhaps footer) on every page, as in the following example:

24 Introduction

[...]

The Nature of Exponential Growth 25

As this only contains information regarding the formatting of the original book (page number and chapter/book title), it should be removed from the corpus text.

Often, we can see that these headers follow the same patterns: the even numbered pages begin with a one-, two- or three-digit number, followed by a series of blank spaces, followed by the title of the chapter. They also always begin on a new line and end with a line break.

The odd-numbered pages, on the other hand, include the chapter title, followed by a series of blank spaces, followed by a one-, two- or three-digit number.

So, to remove these, we can simply type the following two regexes in the Search box and leave the Replace box blank (replace with nothing = delete).

For the even-numbered pages:

```
\n[0-9]{1,3}(\s)+(Introduction)\s[0-9]\n
```

And for the odd-numbered pages:

```
\n(The Nature of Exponential Growth)(\s)+[0-9]{1,3}\n
```

etc.

- *Removing/modifying existing xml tags:*

Some texts in the corpus may have been obtained in a format that already contains xml markup, which needs to be removed before the proper OMC markup can be inserted. The task of removing and/or modifying existing tags is made much easier with regular expressions.

If one encounters, for instance, <div> tags that mark the beginning and end of a chapter of a work, they can either be removed, or, if appropriate, they can usefully be converted into <section> tags by searching `<div type="chapter" n="([0-9]{1,3})" .*?>` and replacing this with `<section id="$1">`.

Notes and tips

To enable or disable regular expressions in *Visual Studio Code*, use the shortcut **Alt+R**.

If a regular expression is not working in the way that you expect, it is often helpful to break it down into smaller and less complicated expressions, and to test these individually, in order to help identify the problem.

It is best to use the period or dot character sparingly. Often, a character class or negated character class is faster and more precise.