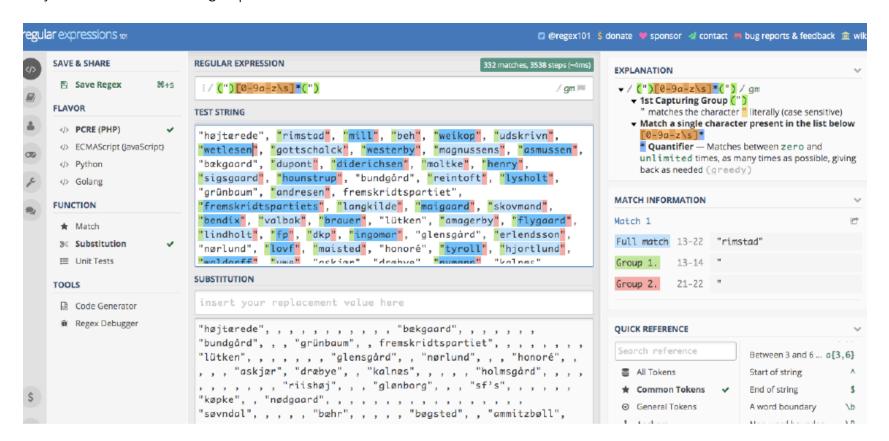# W1: Regular expressions

Here are the links to my two solved exercises:

1: https://regex101.com/r/p5096s/1
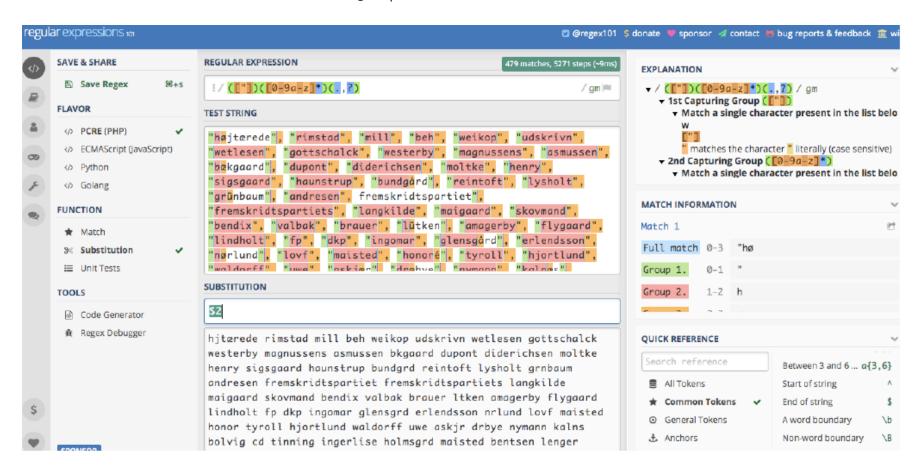
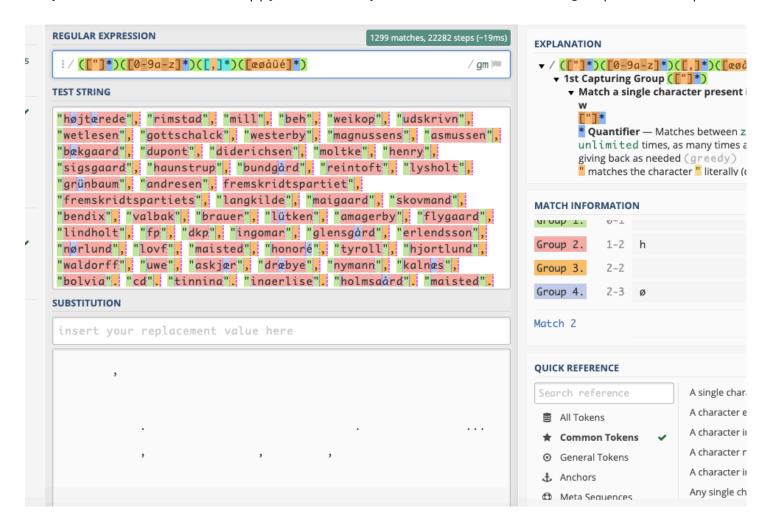2: https://regex101.com/r/COHXGB/1

## EXERCISE 1

This was the first time for me working with regular expressions, so I tried a lot of things to make it work. I wanted to isolate the "", so I tried to do that, but they went into two different groups
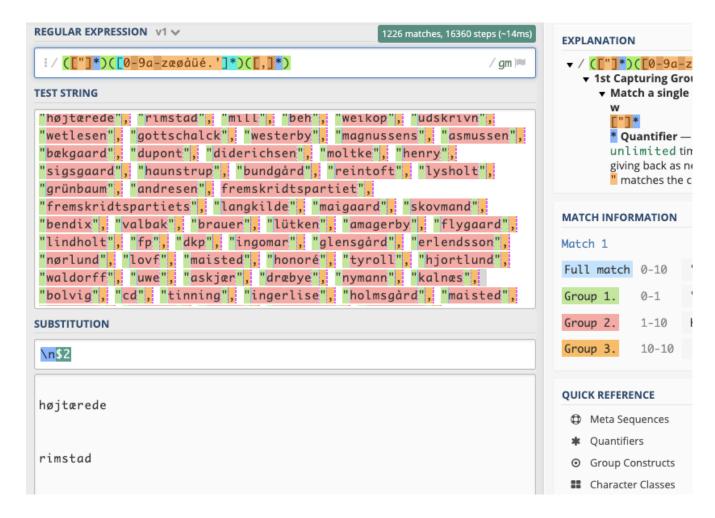
I added [], and then the "" went into the same group. Still, some words weren't marked

I could see that something was wrong in the words who had characters like "æøåüé". Therefore, I included these in my regular expression. They were marked and I was happy, but still, they were marked as their own group instead of part of the word
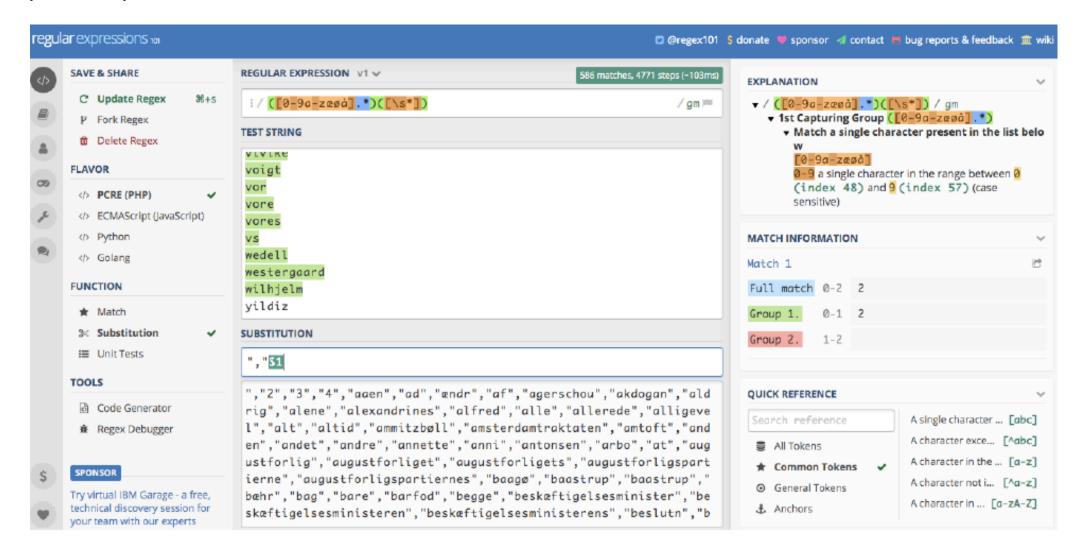
I tried to put these "æøåüé" a different place in the regular expression, in extension of the words, and suddenly it worked! Then I figured out that to finish my work in the substitution box, I needed to use "/n". Then my work was finished. But I did not manage to sort the words alphabetically, is that possible? And sure, there might be smarter ways to do it… Probably a way where you do not have to write all the things that you want it to mark, but rather: mark what's in between ""

**REGULAR EXPRESSION** v1 ⌄                                    1226 matches, 16360 steps (~14ms)

⫶ / `(["]*)([0-9a-zæøåüé.']*)([,]*)`                                            / gm ⚑

**TEST STRING**

```
"højtærede", "rimstad", "mill", "beh", "weikop", "udskrivn",
"wetlesen", "gottschalck", "westerby", "magnussens", "asmussen",
"bækgaard", "dupont", "diderichsen", "moltke", "henry",
"sigsgaard", "haunstrup", "bundgård", "reintoft", "lysholt",
"grünbaum", "andresen", fremskridtspartiet",
"fremskridtspartiets", "langkilde", "maigaard", "skovmand",
"bendix", "valbak", "brauer", "lütken", "amagerby", "flygaard",
"lindholt", "fp", "dkp", "ingomar", "glensgård", "erlendsson",
"nørlund", "lovf", "maisted", "honoré", "tyroll", "hjortlund",
"waldorff", "uwe", "askjær", "dræbye", "nymann", "kalnæs",
"bolvig", "cd", "tinning", "ingerlise", "holmsgård", "maisted",
```

**SUBSTITUTION**

`\n$2`

højtærede


rimstad

**EXPLANATION**

⌄ / `(["]*)([0-9a-z`
  ⌄ **1st Capturing Gro**
    ⌄ **Match a single**
       w
       `["]*`
       **\*** Quantifier —
       unlimited tim
       giving back as n
       `"` matches the c

**MATCH INFORMATION**

Match 1

| Full match | 0-10 |
|---|---|
| Group 1. | 0-1 |
| Group 2. | 1-10 |
| Group 3. | 10-10 |

**QUICK REFERENCE**

⊕  Meta Sequences

✳  Quantifiers
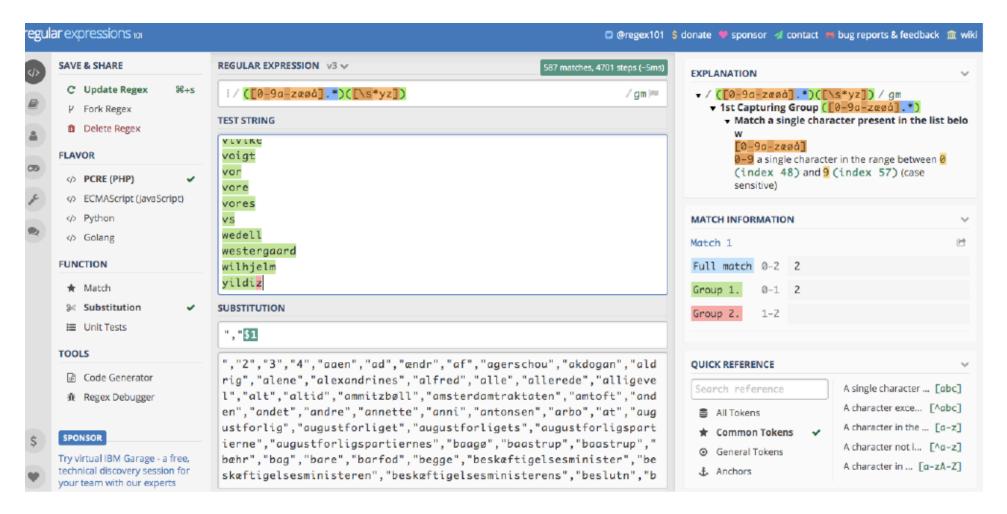
⊙  Group Constructs

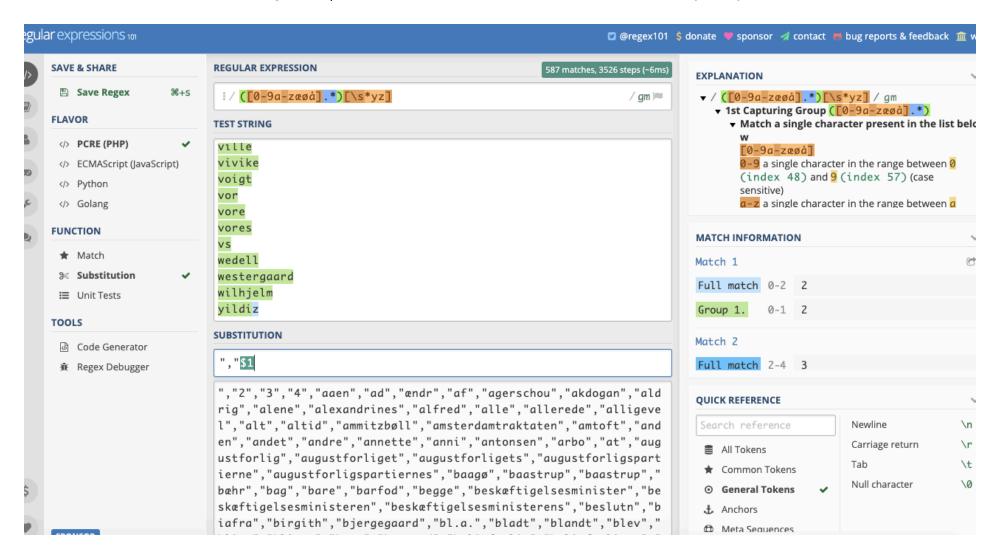▦  Character Classes

## EXERCISE 2

I tried to do the same as in exercise 1, and quickly all the words were marked. I also added "([\s]) to remove the spaces. In the substitution box, I added a comma so it would look like a stopword list. One thing I did not understand was why the world "yildiz" was not part of the list. Do you know why?

Anyway, I tried different things to include it, but it could probably have been done smarter. Here it still did not mark yildiz as one group

After I deleted the "()" in mt last regular expression, it worked. But there must be an easier way. Do you know that?



As I am completely new to this and I mostly did this exercise by just trying different things, I have a few questions…

## QUESTIONS

Related to exercise one:
- Is there a smarter way to just nominate all the characters in between "", regardless of what they are?
- Is there a way to list the words alphabetically?

Related to exercise two:
- Why was the word yildiz not marked?