

Sprachverarbeitung: Übung

SoSe 24

Janis Pagel

Department for Digital Humanities, University of Cologne

2024-04-25

Please submit your solution via Ilias, either as a Jupyter Notebook (.ipynb, you can export your Notebook in Jupyter by going to File > Download) or as a Python script (.py) if you are not working in Jupyter.

Exercise 1.

Write a single regular expression in Python that matches all the strings in the following lists, and only these strings:

- `["pit", "spot", "spate", "slap two", "respite"]`
- `["rap them", "tapeth", "apth", "wrap/try", "sap tray", "87ap9th", "apothecary"]`
- `["affgfkking", "rafgkahe", "bafghk", "baffgkit", "affgfkking", "rafgkahe", "bafghk", "baffg kit"]`
- `["www.google.com", "http://www.google.pl/search?q=exercise+programming", "https://www.google.cz/search?q=exercise&var=variable"]`
- `["+424 161 727 363", "+000 000 000 000", "123 321 765", "654 789 123"]`

Try to make your regular expressions as compact as possible.

These exercises are taken from <https://regex.sketchengine.co.uk/>, you can check their website if you want more exercises to practice on.

Exercise 2.

On <https://lehre.idh.uni-koeln.de/site/assets/files/5151/verwandlung.txt> you find the full text of Franz Kafka's *Die Verwandlung*. Read in the text with Python and write regular expressions that find the following:

- Find all chapter numbers (and only the chapter numbers)

- Find all words followed by either a question mark or an exclamation mark (including the question mark or exclamation mark)
- Find all occurrences of text in-between two commas within a sentence. For example, your regex should find “, wenn er den Kopf ein wenig hob,”, but not “, fand er sich in seinem Bett zu einem ungeheueren Ungeziefer verwandelt. Er lag auf seinem panzerartig harten Rücken und sah,”
- Find all occurrences of direct speech
- Find all nominalizations that start with either “Be” or “Er” and end with “ung”. Do you also find nouns that are not a nominalization and if yes, which ones? Would it be possible to exclude these “false” hits with regular expressions in a general way?