

Aufgabe 1: Alice im Wunderland

Team-ID: 01093

Team-Name: Kimi and his Droogs

Bearbeiter/-innen dieser Aufgabe:
Kimi Müller

20. November 2022

Inhaltsverzeichnis

1	Lösungsidee	1
2	Umsetzung	1
3	Beispiele	1
4	Quellcode	2

Anleitung: Trage oben in den Zeilen 8 bis 11 die Aufgabennummer, die Team-ID, den Team-Namen und alle Bearbeiter/-innen dieser Aufgabe mit Vor- und Nachnamen ein. Vergiss nicht, auch den Aufgabennamen anzupassen (statt „ \LaTeX -Dokument“)!

Dann kannst du dieses Dokument mit deiner \LaTeX -Umgebung übersetzen.

Die Texte, die hier bereits stehen, geben ein paar Hinweise zur Einsendung. Du solltest sie aber in deiner Einsendung wieder entfernen!

1 Lösungsidee

Abstrahiert ist das Lückenwort eine Buchstabenfolge beliebiger Länge ≥ 1 . Daher suchen wir einfach nach dem Kontext des Lückenwortes, die Lücke wird da mit dieser Buchstabenfolge ersetzt.

2 Umsetzung

Mithilfe des regulären Ausdrucks " w+ " für das Lückenwort suchen wir nach Übereinstimmungen. " w " ist eine Abkürzung für alle Wortbestandteile, so muss sich keine Sorge um Umlaute o.ä. gemacht werden. " + " steht für "1 oder mehrere Instanzen des vorherigen Worts"

3 Beispiele

Die inkludierten Beispiele:

```
1 ['das kommt mir gar nicht richtig vor']
2 ['ich mu doch clara sein', 'ich mu in clara verwandelt']
3 ['fressen spatzen gern katzen', 'fressen katzen gern spatzen']
4 ['das spiel fing an', 'das publikum fing an']
5 ['ein sehr schöner tag']
6 ['wollen sie so gut sein']
```

Zeile 3 sticht als besonders hervor, wo zum ersten Mal unerwartetes Verhalten auftrat. Wird das Buch unabhängig von Zeilenabsätzen durchsucht, kommt 'fressen katzen gern spatzen' doppelt vor. Dem wird vorgebeugt, indem identische Treffer gelöscht.

```
1 ['alice in späterer zeit']
2 ['aufschrift eingemachte']
```

Beispiel (1) ist relevant, da im Buch ein Teil des gesuchten innerhalb einer neuen Zeile ist. Daraus ergibt sich sowohl das Problem, dass nicht nach Leerzeichen gematcht werden muss, sondern alle "leeren Zeichen"(neben Leerzeichen auch tabs, neue Zeile/carriage return etc.) gematcht werden können, als auch dass die neue Zeile im Ergebnis übernommen werden würde. Dies wurde gelöst, indem alle "leeren Zeichen" zu Leerzeichen konvertiert wurden.

Beispiel (2) wurde erstellt, da es Sonderzeichen enthielt.

4 Quellcode

Der Vollständigkeit halber möchte ich hier den Proof-of-concept in bash zeigen, der mit den genannten Sonderfällen nicht funktioniert:

```
1 gaptext='cat $1';
2 cat $2 | grep -iEo "${gaptext//_/"\w+"}"
```

Zur eigentlichen Lösung in Python:

```
1 def convert_gap_text_to_regex(gap_text):
2     """
3     Converts a gap text to a precompiled regular expression
4     """
5     re_gap_text = gap_text.strip()
6     re_gap_text = re_gap_text.replace(" ", "\ ")
7     re_gap_text = re_gap_text.replace("_", "\w+")
8     return re.compile(re_gap_text)
9
10 def tidy_up_source_book(source_book):
11     """
12     Eliminates special Characters and converts any whitespace to a space char
13     """
14     clean_source = re.sub("[,.:!?!;¿ñ_]*", "", source_book)
15     clean_source = re.sub("\s+", " ", clean_source)
16     clean_source = clean_source.lower()
17     return clean_source
18
19
20 def find_all_results(interference, favourite_book):
21     """
22     This function mends the other two into a function that outputs all matches.
23     """
24     regex = convert_gap_text_to_regex(interference)
25     book = tidy_up_source_book(favourite_book)
26     return list(set(regex.findall( book)))
```