

Computational Methods for Document Analysis

02

Mirco Dietrich
Maurice-Roman Isele

25. April 2016

Exercise 2

Problem 1:

People don't always use a white-space after a punctuation mark or use multiple marks. Replacing each punctuation mark by an empty string will glue the last and first word of sentences together. So we needed to replace each punctuation mark by a whitespace and then reduce multiple whitespaces to one.

Problem 2:

' Should be removed around words if it is used as a quotation mark, but shouldn't be removed when it is used as in a abbreviation(e.g. "*you're*"). Since there are not that many quotes with ' we decided to not remove it. There would be a solution to remove the' from quotes although it would be a bit of a hack:

You could read in a file where the `String[]` with the tokens was saved with `Arrays.toString`. Then replace ' (which is always at the end of a quote), " (which is always at the beginning of a quote) and all other no-characters with whitespaces, split at whitespaces and save the new tokens in the Review object.

Problem 3:

Removing ', ' and ' / ' splits numbers separated by them(e.g. decimals, fractions). We don't have a solution for that.

Exercise 3

One of the obvious problems was to identify the abbreviations used by the people reviewing. We solved this by looking at the data and furthermore considering more common abbreviations used in the English language. Additionally, any name abbreviations are possible, e.g. "*Mr. K.*", "*T.J.*", and so on.

