

## Report

During the third step of part A, we extend our program to implement a technique called smoothing.

### Function

For natural language purposes - *regarding calculating probabilities* - smoothing is used to make a more realistic view of the occurrences of sentences in a language. Without smoothing, the values of the calculated probabilities will result in a skewed view on possible sentences. N-grams that are assigned a value 0 - *because they do not occur in the corpus* - will be regarded as an impossible combinations of words, rather than an unlikely combination of words. Due to the distinct computational ceiling it is impossible to calculate all exact probabilities of n-grams. To solve this problem, we apply smoothing to be able to make more realistic assumptions about these probabilities. This results in better generalizations and thus we obtain a better view of the language and the occurring n-grams, even though the corpus we used is finite.

### General

The fundament of smoothing methods is commonly defined as '*steal from the rich and give to the poor*'. B

### Methods

1. **Add-1 smoothing (Laplace smoothing).** This is the most trivial type of smoothing. Every occurrences count is increased by one. Then, all frequency counts are normalized. This is necessary because the resulting probabilities must be of the '*normal distribution*'-type.
2. **Good-Turing smoothing.** A different, more sophisticated method to calculate probabilities of unseen events.  $r^*$  is considered the new assigned probability and is calculated through the following equation.

$$r^* = \frac{(r + 1) \frac{n_{r+1}}{n_r} - r \frac{(k+1)n_{k+1}}{n_1}}{1 - \frac{(k+1)n_{k+1}}{n_1}}$$

We only applied this smoothing for bigrams that appeared less than 5 times in the corpus.

## Answers

The percentages of sentences that are assigned a probability of 0 are:

Add-1 smoothing	
Good-Turing smoothing	

The first five sentences that are assigned a probability of 0 are:

	<i>Add-1 smoothing</i>	<i>Good-Turing smoothing</i>
1	sen1	sen1
2	sen2	sen2
3	sen3	sen3
4	sen4	sen4
5	sen5	sen5

# 1 Running

The code can be executed by running the following code.

```
python BerkerGerritseMooijen3.py etc[] etc[] etc[]
```