

Programming and Classification:

3. Simple similarity of texts

Marek Klonowski

Suggested deadline: 15.04.2019

You will need NLTK <https://www.nltk.org/>.

21. ★ For a given bitstring **b** list all bitstrings **b'**, such that the Hamming distance between **b** and **b'** is equal 1.
22. ★ Construct a function that returns a Jaccard similarity for two sets. Beware that this function needs to check if at least one of the sets is nonempty.
23. ★ Construct a function that computes Jaccard similarity for two strings treated as bags of words.
24. ★★ (use NLTK) List all words in `text1` with edit distance from the word `dog` smaller than 4. Hint: you can safely reject all long words without computations (why?).
25. ★★ (use NLTK) Let `text1` - `text9` be bags of words. Compute similarity between all pairs of texts.
26. ★★ (use NLTK) Let us consider a metric space (S, d) , where S is the set of words from `text1` and d is the Hamming distance. Find diameter of (S, d) .
27. ★★★ (use NLTK) Construct a dictionary that assigns each pair of consecutive words in `text1` the Jaccard similarity between them.
28. ★★★ (use NLTK). For two words v and w , let *relative edit distance* be the Levenshtein distance between v and w divided by the sum of lengths v and w . Find two **different** words in `text2` with minimal relative edit distance.
29. ★★★ For a given bitstring **b** and a natural number r list all bitstrings **b'**, such that the Hamming distance between **b** and **b'** is equal n .
30. ★★★ Construct a function that for a given string and a natural number k returns a **set** of all its k -shingles.