

NPFL097 Assignment 3

Word Clustering and Component Analysis

David Mareček

December 2021

In this assignment, you will get a vocabulary of English words with corresponding 512-dimensional vector representations. These word vectors are used for predicting the next English word in a generated English sentence and were obtained from the Czech-English Neural Machine Translation model.

Your task is to analyze this highly-dimensional vector space and find the main features by which the words are distributed. The preferred programming language is Python.

Send me your final source code and commented graphs and results by email: marecek@ufal.mff.cuni.cz.

Data

Download the file `nmt-en-dec-512.txt` containing word vectors of the 5,000 most frequent words from English fiction. There is one word on each line followed by its vector representation (512 tab-separated real numbers).

Tasks

1. Visualize the word vectors in 2D using the T-SNE method. Use the `sklearn.manifold.TSNE` function. Learn how T-SNE works and what parameters it has. Draw a scatterplot with a reasonable amount of labels, so that it is readable and you can identify different clusters of words. (1 pt)
2. Cluster the word vectors in the original 512-dimensional vector space using the methods of K-means, GaussianMixture, and Agglomerative clustering. Use the functions implemented in the `sklearn` library. Try different number of clusters (3,5,10) and different linkage strategies (ward, single, complete). Visualize the clusters using different colors in T-SNE. Compute the Silhouette coefficients of the individual clusterings. Use

