

This thesis focuses on automatic searching for clusters of topically similar texts in large text collection. We introduce an algorithm for finding the clusters and a method of optimizing its parameters using machine learning techniques. The algorithm is implemented and experimentally evaluated. For evaluation we use a manually annotated collection of Czech documents, which contains a set of sample clusters chosen and tagged by a human annotator, and a huge collection of newspaper articles. Experiments show that the output of our algorithm fulfills our expectation and gives clusters of topically similar texts.