Collections of texts are considered as a valuable source of information for applied economic analysis. Recent developments in the access to large sets of documents, e.g., scientific abstracts, articles, news items, social media messages or statements of different institutions, and in the methods developed for extracting information from texts increase the interest in this type of data. However, the knowledge about the performance of these methods, in particular when combined with the usual econometric methods is still rather limited. Therefore, the objective of the TEXTMOD project is to contribute to the development of methods and to improve the understanding of how the information obtained from text mining can be incorporated in econometric models. Thereby, the focus is on multivariate time series models. The indicators are constructed using models which try to identify relevant themes in large collections of documents without human intervention. An example of text-based time series, which can be of interest in economic research and can add information content to classical real economic indicators, is a topic trend describing how the importance of a given topic (e.g. related to inflation) changed over time. While a substantial number of methods have been proposed over the last few years for identifying topics and their trends over time, there is little evidence on the statistical properties of these procedures, their relative performance and their interaction with more traditional modelling approaches. Consequently, a central aim of the project is to investigate sensitivity to parameter settings, robustness to variations of the textual sample and uncertainty associated with these algorithms. In the project, additional methods for comparing the results of topic modelling across samples or resulting from different methods will be proposed. In a further important step, different methods for deriving trends in topics will be considered and finally the consequences of including them in time series models, e.g., the widely used vector autoregressive model, will be studied. Special emphasis will be put on the appropriate interpretation of results, evaluation of additional insights from using text-based data and rigorous measurement of the estimation uncertainty which will be captured by means of joint confidence bands. The methods will be applied to study the relationships between real economic indicators and trends in topics found for scientific corpora in economics from Poland and Germany.