# Convergence of the Algorithm of Additive Regularization of Topic Models

## I. A. Irkhin[1,*] and K. V. Vorontsov[1,**]

**Abstract**—The problem of probabilistic topic modeling is as follows. Given a collection of text documents, find the conditional distribution over topics for each document and the conditional distribution over words (or terms) for each topic. Log-likelihood maximization is used to solve this problem. The problem generally has an infinite set of solutions and is ill-posed according to Hadamard. In the framework of Additive Regularization of Topic Models (ARTM), a weighted sum of regularization criteria is added to the main log-likelihood criterion. The numerical method for solving this optimization problem is a kind of an iterative EM-algorithm written in a general form for an arbitrary smooth regularizer as well as for a linear combination of smooth regularizers. This paper studies the problem of convergence of the EM iterative process. Sufficient conditions are obtained for the convergence to a stationary point of the regularized log-likelihood. The constraints imposed on the regularizer are not too restrictive. We give their interpretations from the point of view of the practical implementation of the algorithm. A modification of the algorithm is proposed that improves the convergence without additional time and memory costs. Experiments on a news text collection have shown that our modification both accelerates the convergence and improves the value of the criterion to be optimized.

**Keywords:** natural language processing, probabilistic topic modeling, probabilistic latent semantic analysis (PLSA), latent Dirichlet allocation (LDA), additive regularization of topic models (ARTM), EM-algorithm, sufficient conditions for convergence.

[1]Moscow Institute of Physics and Technology (National Research University), Dolgoprudnyi, 141701 Russia
e-mail: *ilirhin@gmail.com, **k.v.vorontsov@phystech.edu