# Università degli Studi di Udine

### Dipartimento di Scienze Matematiche, Informatiche e Fisiche

# Alpen-Adria-Universität Klagenfurt

### Faculty of Technical Sciences

### Corso di Laurea Magistrale in Informatica Internazionale
### Joint Degree in International Informatics

Master Thesis

# Improving the Efficiency of Retrieval Effectiveness Evaluation: Finding a Few Good Topics with Clustering and Dimensionality Reduction

*Supervisor:*
Prof. Stefano Mizzaro

*Candidate:*
Kevin Roitero

*Co-supervisor:*
Prof. Klaus Schoeffmann

Academic year 2015-2016

# Abstract

We consider the issue of using fewer topics in the effectiveness evaluation of information retrieval systems. Previous work has shown that using fewer topics is theoretically possible: there exist smaller topic sets that evaluate a population of systems in almost the same way as a full and larger topic set. One of the main issues that remains to be solved is how to find such a small set of a few good topics. To this aim, in this thesis we use a novel approach based on clustering of topics together with dimensionality reduction.

We present various approaches: we consider an a posteriori approach, i.e., we build clusters of topics only after the evaluation exercise (and the relevance assessments) has already been performed; an a priori approach, i.e., before any human assessment has been performed, and thus without using the features produced during assessment; and some approaches that occur during the evaluation process.

We show that clustering is effective in the topic set reduction problem, if the a posteriori approach is considered together with dimensionality reduction (i.e., principal component analysis). We provide furthermore a statistical significance analysis of the obtained result.

We also consider a priori features and features that occur during the evaluation process. The results show that the a priori approaches and the approaches that occur during the evaluation approaches considered are comparable to the random choice of topics. Finally, we provide convincing evidence that straightforward clustering (i.e., without the dimensionality reduction) is not effective in the topic reduction problem.

# Contents