# Analysis and reproducibility of paper "A Word Embedding based Generalized Language Model for Information Retrieval"

D. Ganguly, D. Roy, M. Mitra, G.J.F. Jones

11 february 2020

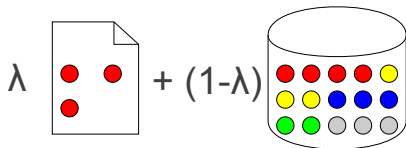Alessio Lazzaron, Matteo Marchiori, Andrea Oriolo, Fabio Piras

# Contents

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

- Lucene... What version?

- Indexing... How?
  - Tokenization
  - Stopwords removal
  - Transposed index

- Base Language Model

$\lambda$ + (1-$\lambda$)

- First approach: let's do it ourselves (fail)

- Second approach: let's try to take example and make it better

- Is it enough? Time constraints...

# References

📄 M. Agosti and G. Silvello, *Slides*, Padua, Italy: academic year 2019-2020.

📄 D. Ganguly, M. Mitra, D. Roy and G. Jones, *A Word Embedding based Generalized Language Model for Information Retrieval*, Dublin, Ireland; Kolkata, India: SIGIR 2015.

📄 D. Ganguly, *GLM GitHub repository*

📄 V. Lavrenko, *Jelinek-Mercer Smoothing video*