# Probabilistic Topic Models
# - TBFY -

Carlos Badenes-Olmedo
Ontology Engineering Group (OEG)
Universidad Politécnica de Madrid (UPM)

✉ cbadenes@fi.upm.es
🐦 @carbadol

oeg-upm.net
github.com/cbadenes

- Each **topic** is a distribution over words
- Each **document** is a mixture of corpus-wide topics
- Each **word** is drawn from one of those topics

*Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet Allocation. Journal of Machine Learning Research, 3(4–5), 993–1022.*

- corpus summarised by topics (ranked list of words)
- scalable representation of documents (mixture of topics)

**Topic 0**

*simplex space*

*documents*

**Topic 1**

**Topic 2**

company's_vision

water_hose

Nando

Otto_Scharmer

complex_transformation_proces

transformation_programs

senior_leadership_team

telecoms_sector

facebook

animal

information

code

string

object

network

class

value

method

emotion

dealer

adapt

human

world

vision

player

story

creation

step

part

transformation

combination

money

• *scalable document repository*

‣ *topics as annotations to **organize**, **summarise** and **search** documents*

‣ ***explore** it in a way that you can index of ideas contained in them*

‣ ***browse** it in a way that you can find documents are about the same kinds of ideas*

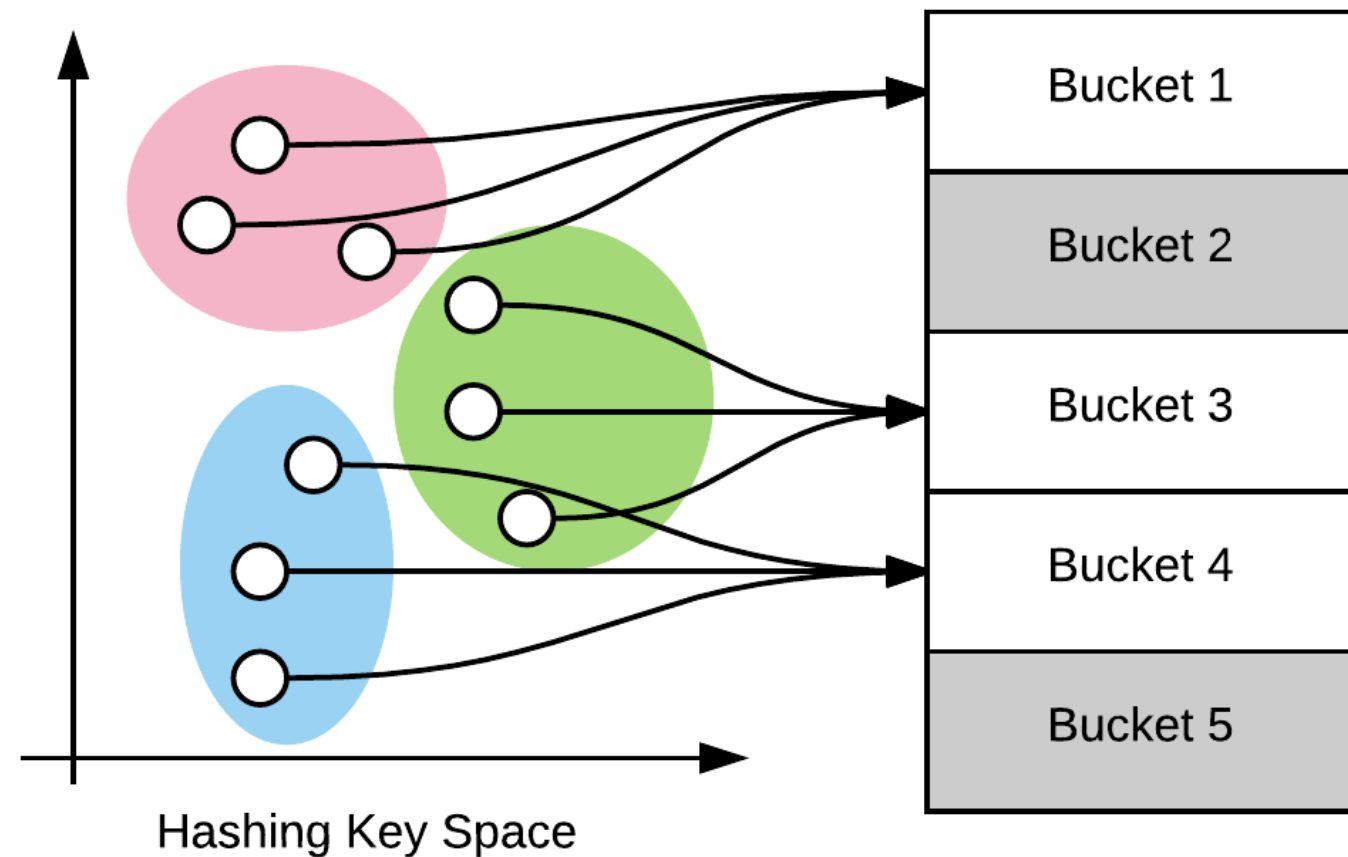- *scalable document repository*
- ***semantic similarity***

 

 

▸ *documents as vectors of topic distributions*

▸ *similarity metric based on **Jensen-Shannon Divergence** (JSD)*

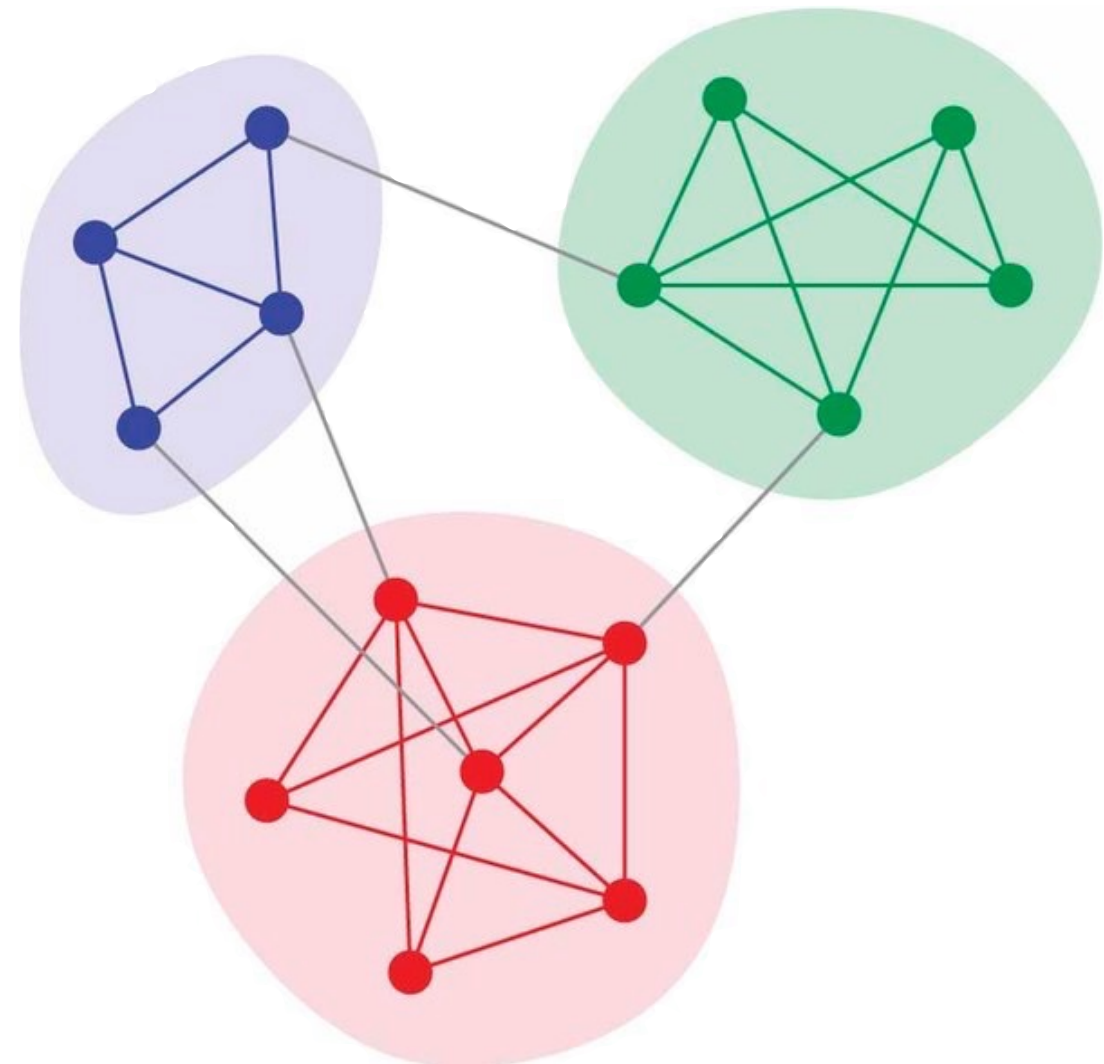▸ *documents and texts*

**0,4953**

- *scalable document repository*
- *semantic similarity*
- ***alerts***

‣ *hashing based on topic distribution*

‣ *non-static approach*

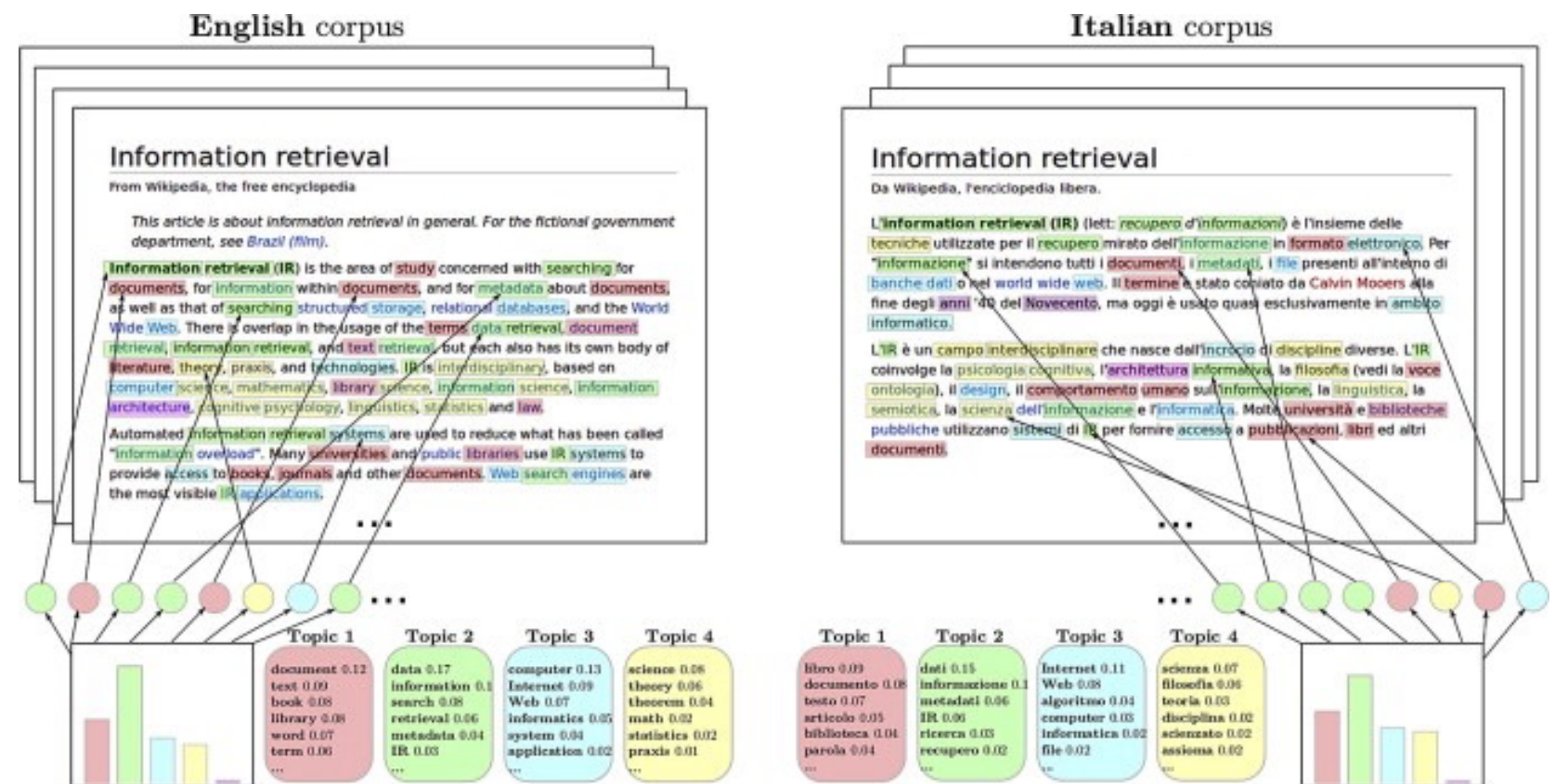‣ *categorical similarities*

‣ *duplicate detection*



Hashing Key Space

- *scalable document repository*
- *semantic similarity*
- *alerts*
- **semantic exploration**

▸ *topics linked to entities from external Knowledge Graph*

▸ *multi-valued classification*

▸ *hierarchical-classification*

▸ *profiling*

- *scalable document repository*
- *semantic similarity*
- *alerts*
- *semantic exploration*
- ***multi-lingual recommendation***

‣ *topics from different languages (models)*

‣ *automatic alignment*



Vulic, Ivan, Wim De Smet, Jie Tang and Marie-Francine Moens. "Probabilistic topic modeling in multilingual settings: An overview of its methodology and applications." *Inf. Process. Manage.* 51 (2015): 111-147.