

Perspectives on Data Access And Use at Scale: Lessons from the Field

Topic Modeling on the Kadi Registers

Dr. Sümeyye Akça

Assistant Professor
Marmara University
Information and Records Management
Turkey/Istanbul
sumeyyesakca@gmail.com
sumeyye.akca@marmara.edu.tr

Kadi Registers?

- One of the most important sources
- Turkish, Arabic, and Persian
- Maintained in book format
- 15th - 20th centuries
- Judgement for every unique event
- Transcription??



Challenges

- 32 letters Persian and Turkish,
- Writing order is from left to right,
- Vowel point!
- Difficult to distinguish words!
- Separation of lines!
- Horizontal plane!
- Changing writing styles!

Problem?

- Need for authority
- Develop community belonging
- Understand present problems!!

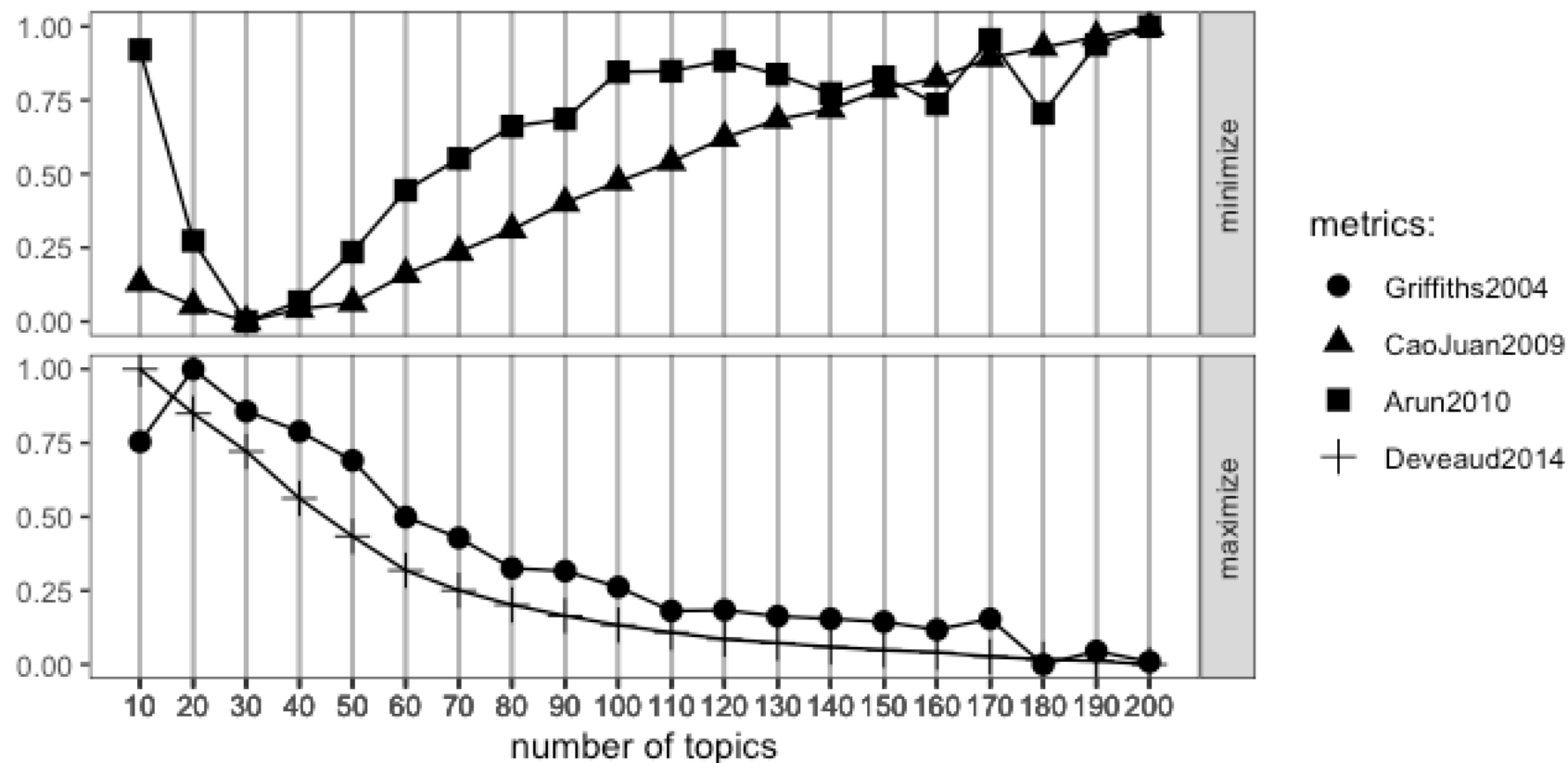
Latent Dirichlet Allocation

- We need algorithmic tools!
- Probability distribution and Topic Modeling!
- LDA algorithm!

The Dataset

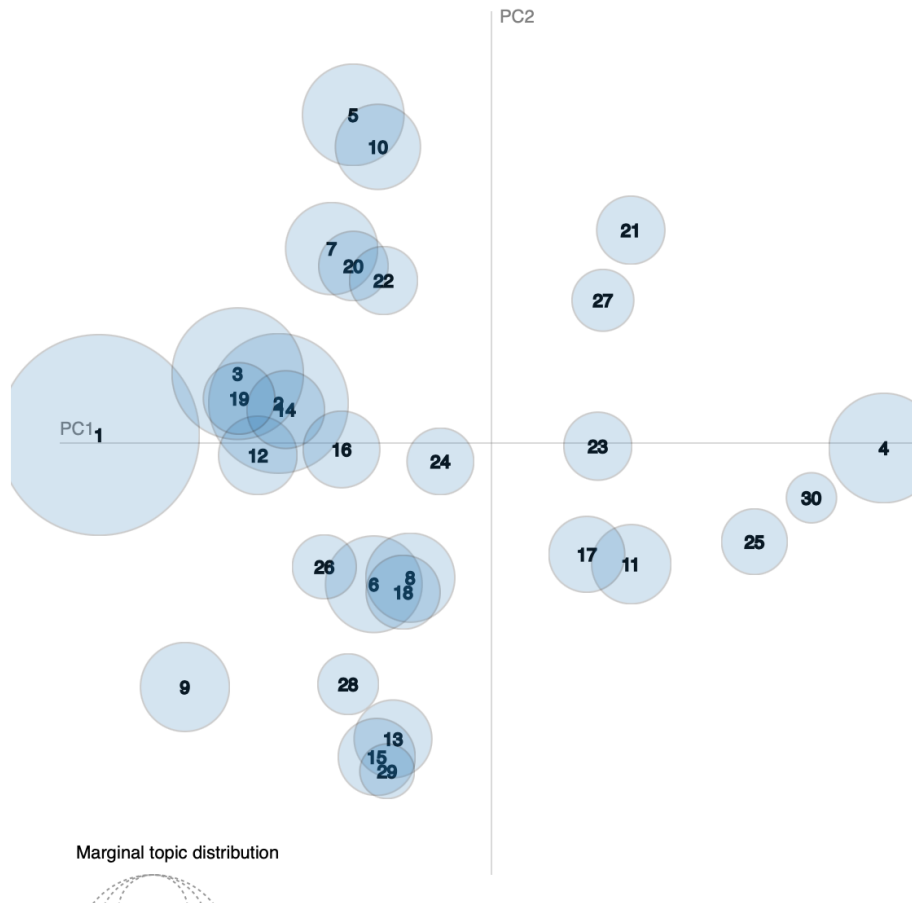
- 23. Kadi Register book
- Consists of 648 single items (problem or case)
- I had transcribed it as my master's thesis
- To determine the topics of the each case =>

Select # of The Topics

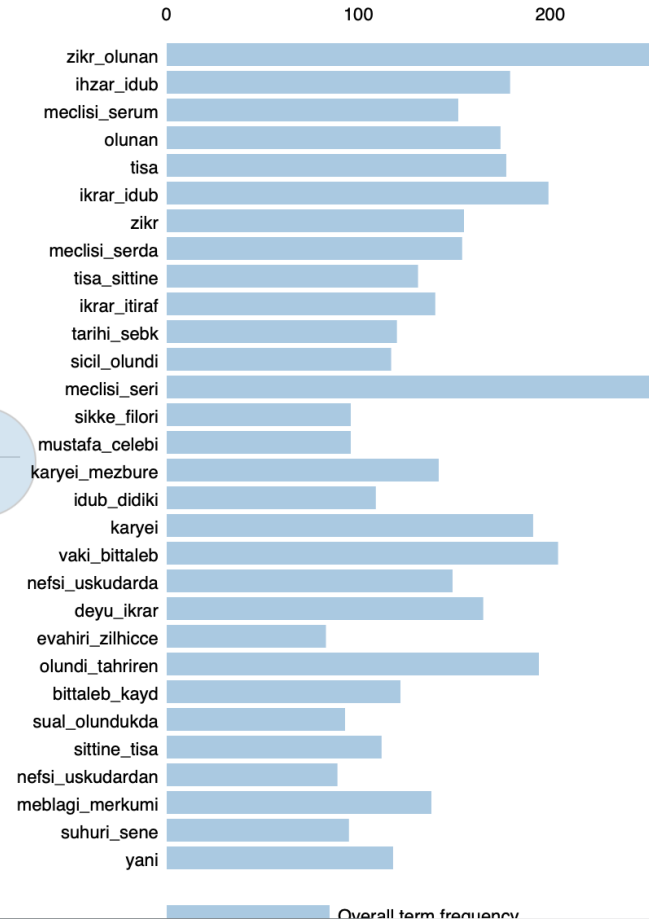


Results

Intertopic Distance Map (via multidimensional scaling)

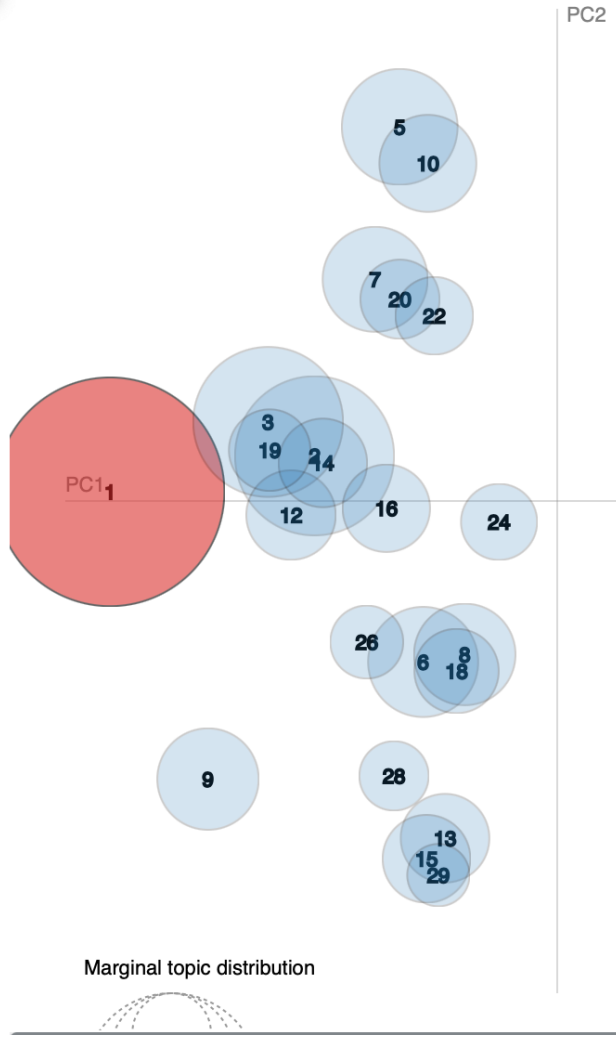


Top-30 Most Salier

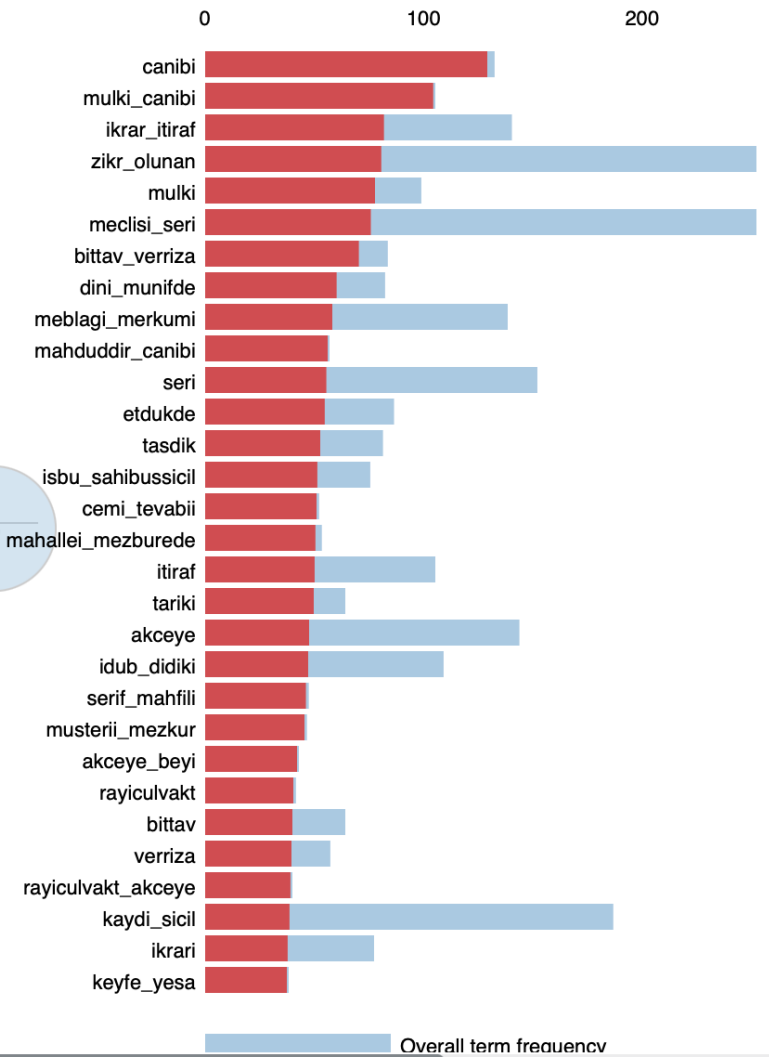


[12]:

Intertopic Distance Map (via multidimensional scaling)

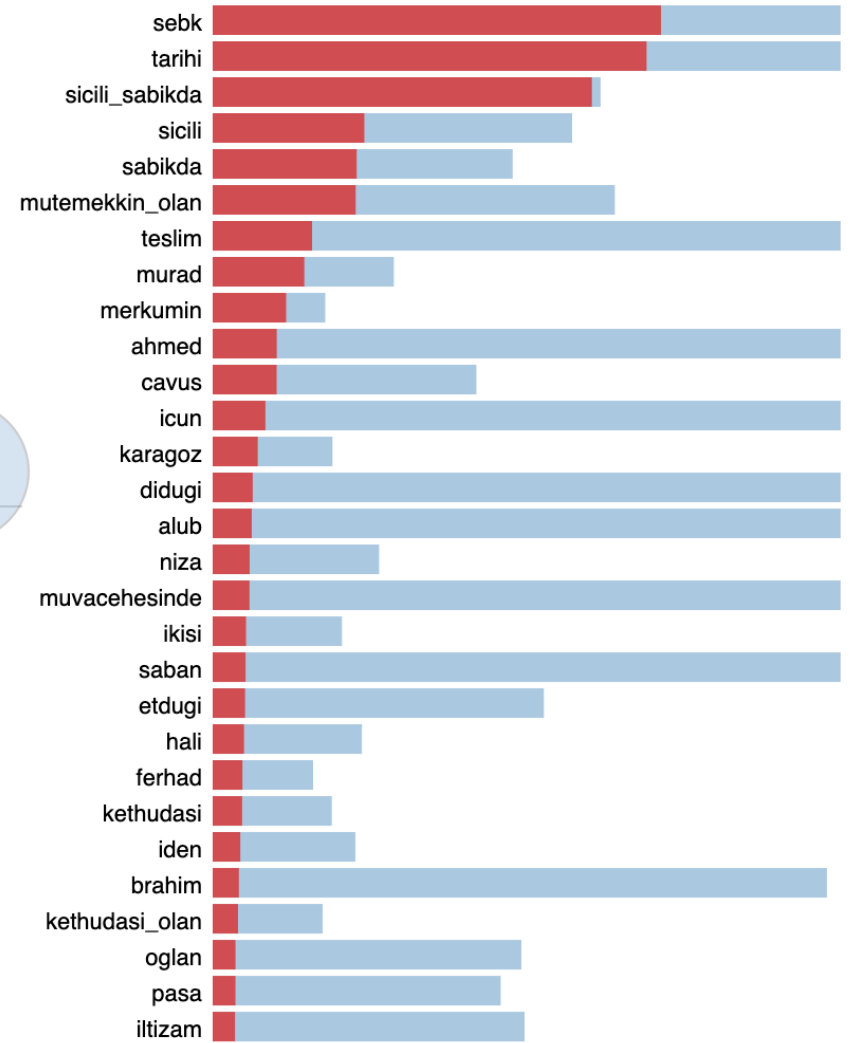
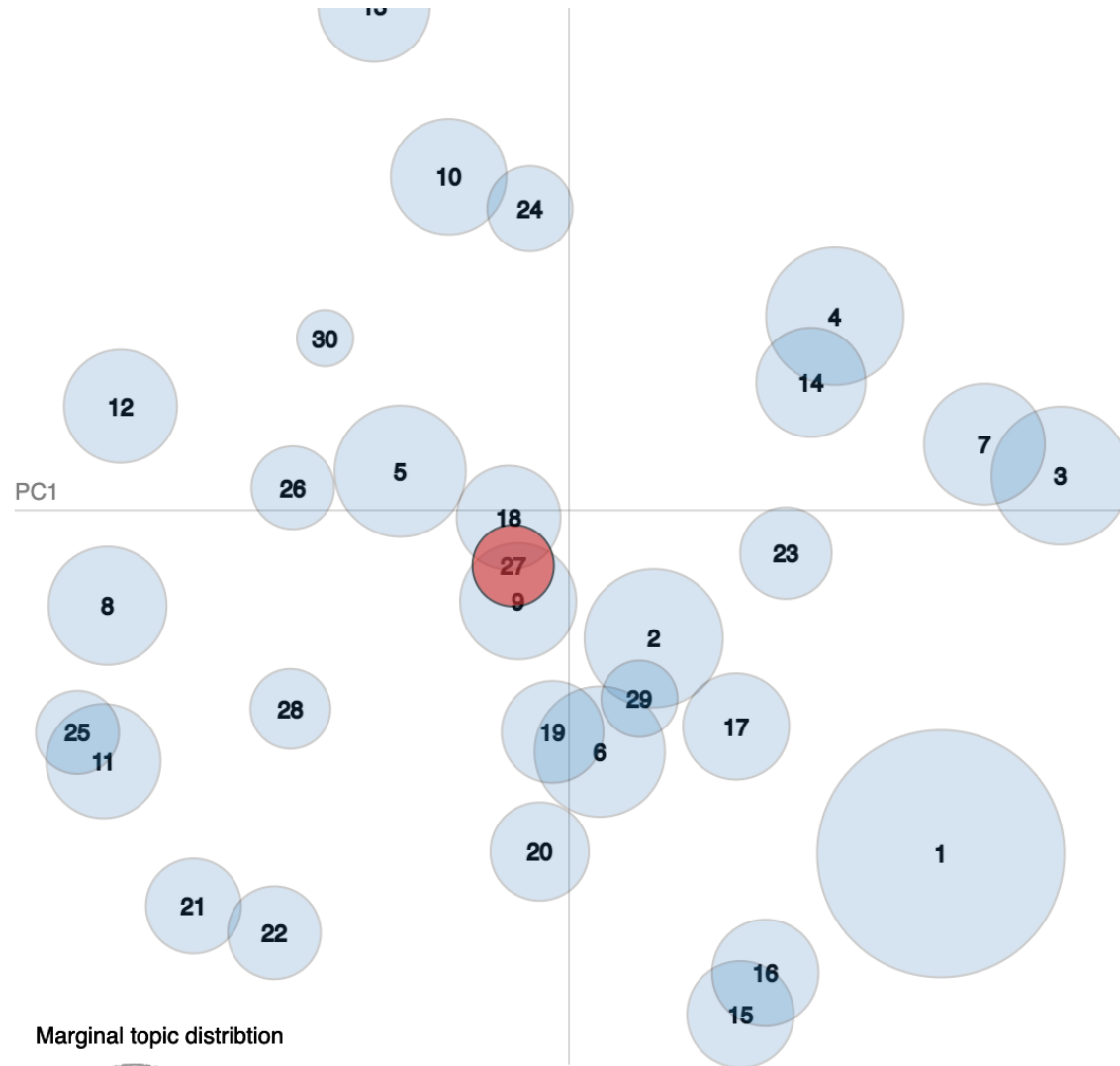


Top-30 Most Relevant Terms for



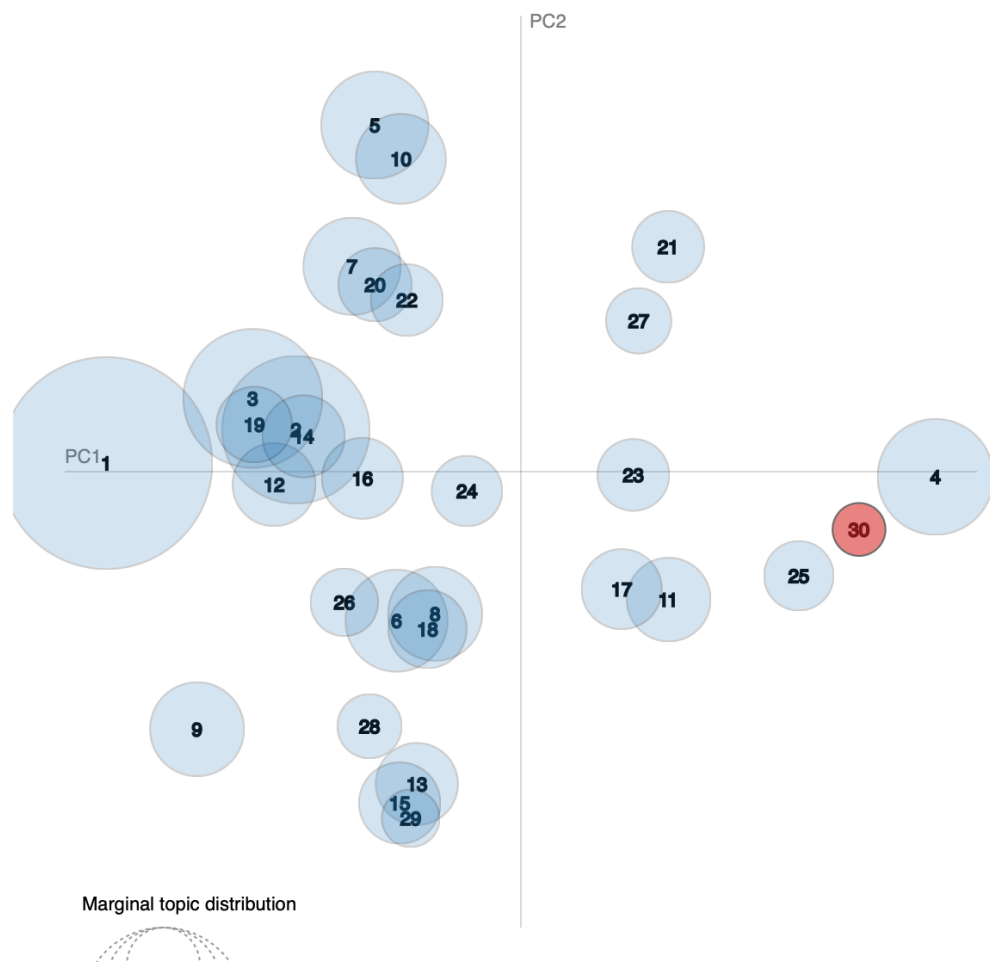
Most Used Words..

- Vechi tahriri huruf oldur ki
- Vechi tahriri kaziyye budur ki
- Ve tarihi ma sebk
- Ve tarihi ma zikru filevvel
- Binti-Bt
- Bin
- Suhud Ma sebk
- Tahriren fi

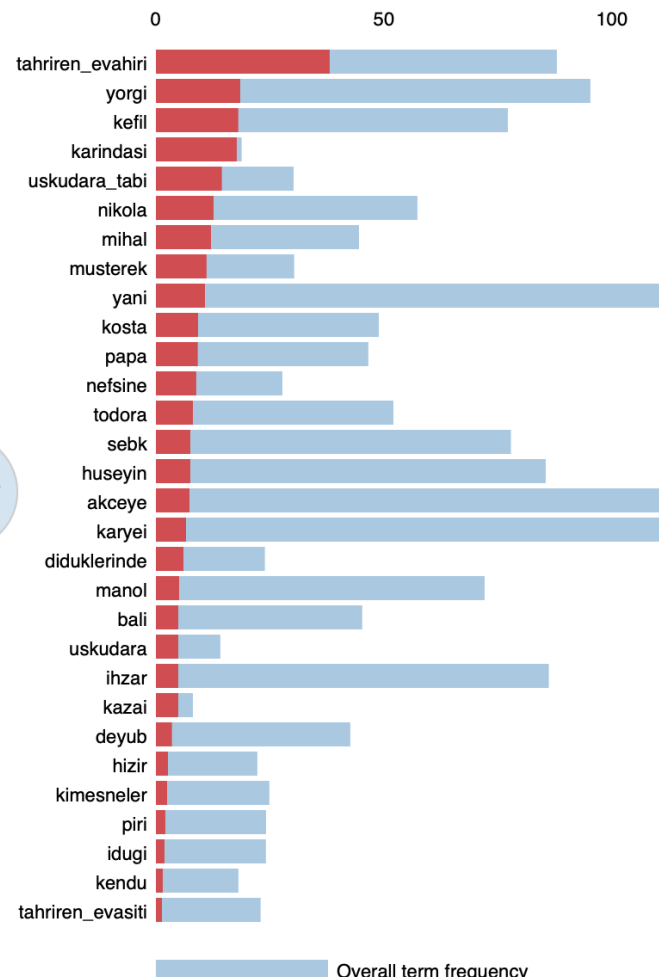


Overall term frequency

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for T_i

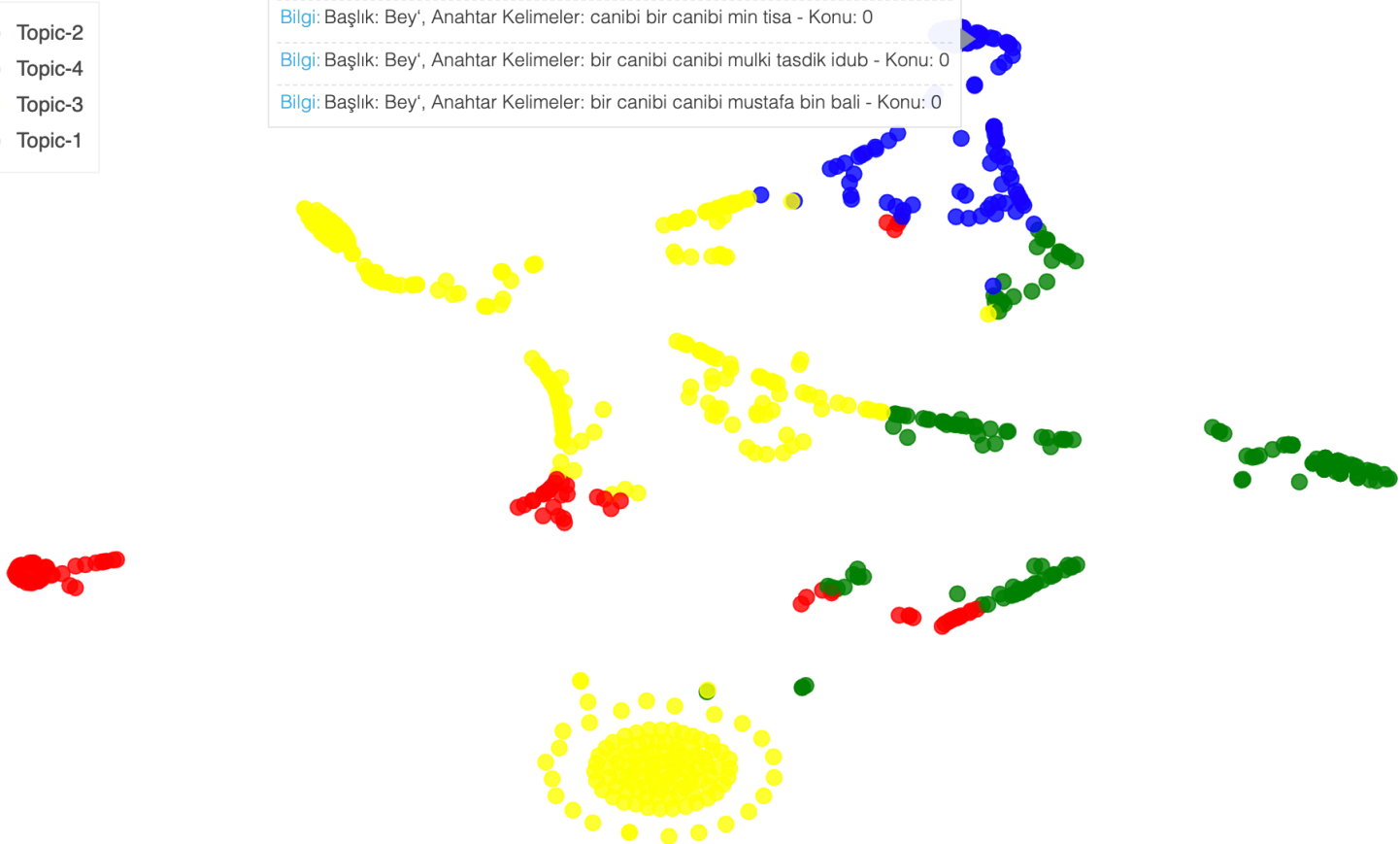


Dimension Reduction

Konuların T-SNE görselleştirmesi

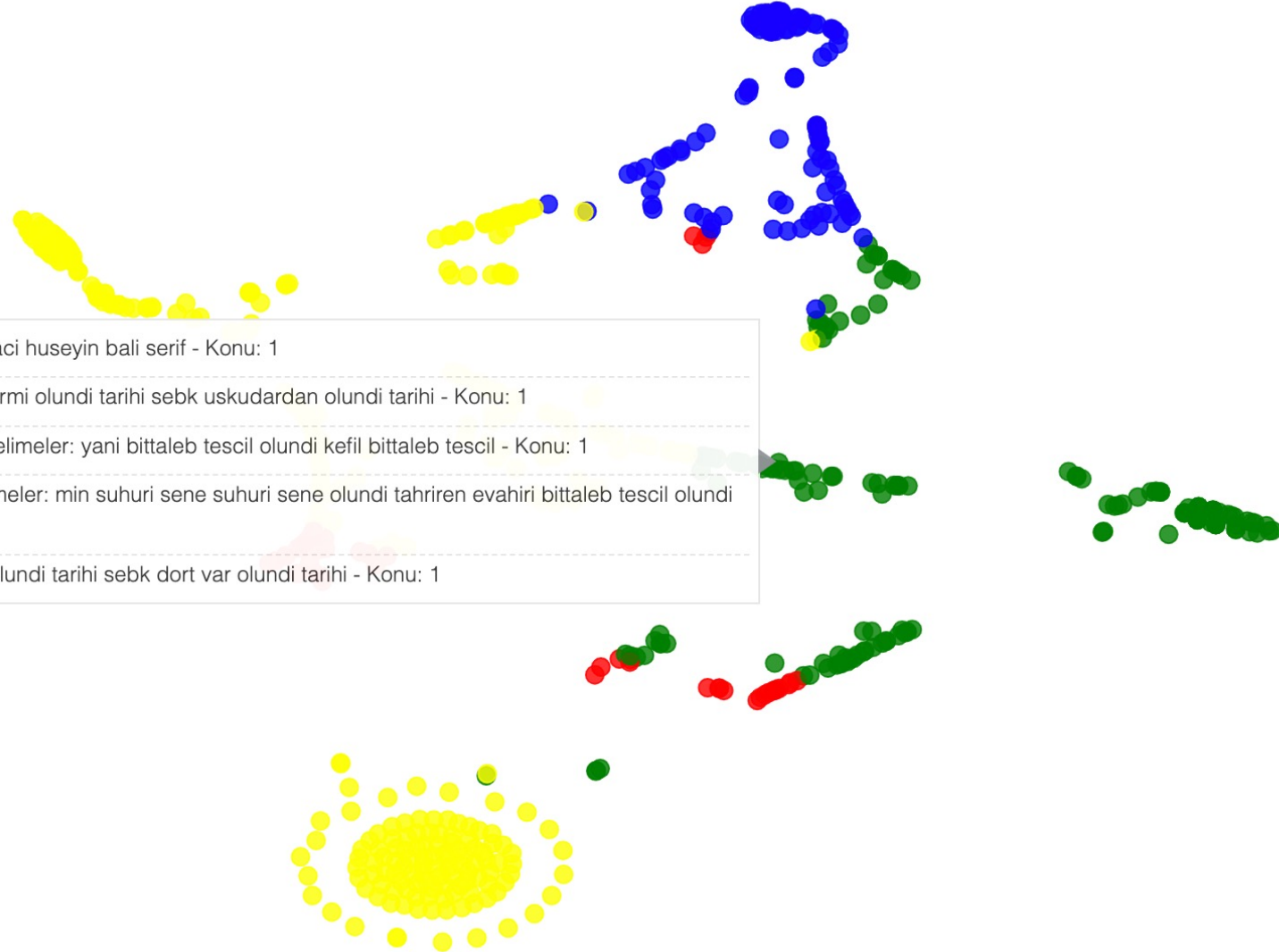
- Topic-2
- Topic-4
- Topic-3
- Topic-1

Bilgi: Başlık: Bey', Anahtar Kelimeler: canibi bir canibi tisa etdim - Konu: 0
Bilgi: Başlık: Bey', Anahtar Kelimeler: canibi bir canibi min tisa - Konu: 0
Bilgi: Başlık: Bey', Anahtar Kelimeler: bir canibi canibi mulki tasdik idub - Konu: 0
Bilgi: Başlık: Bey', Anahtar Kelimeler: bir canibi canibi mustafa bin bali - Konu: 0

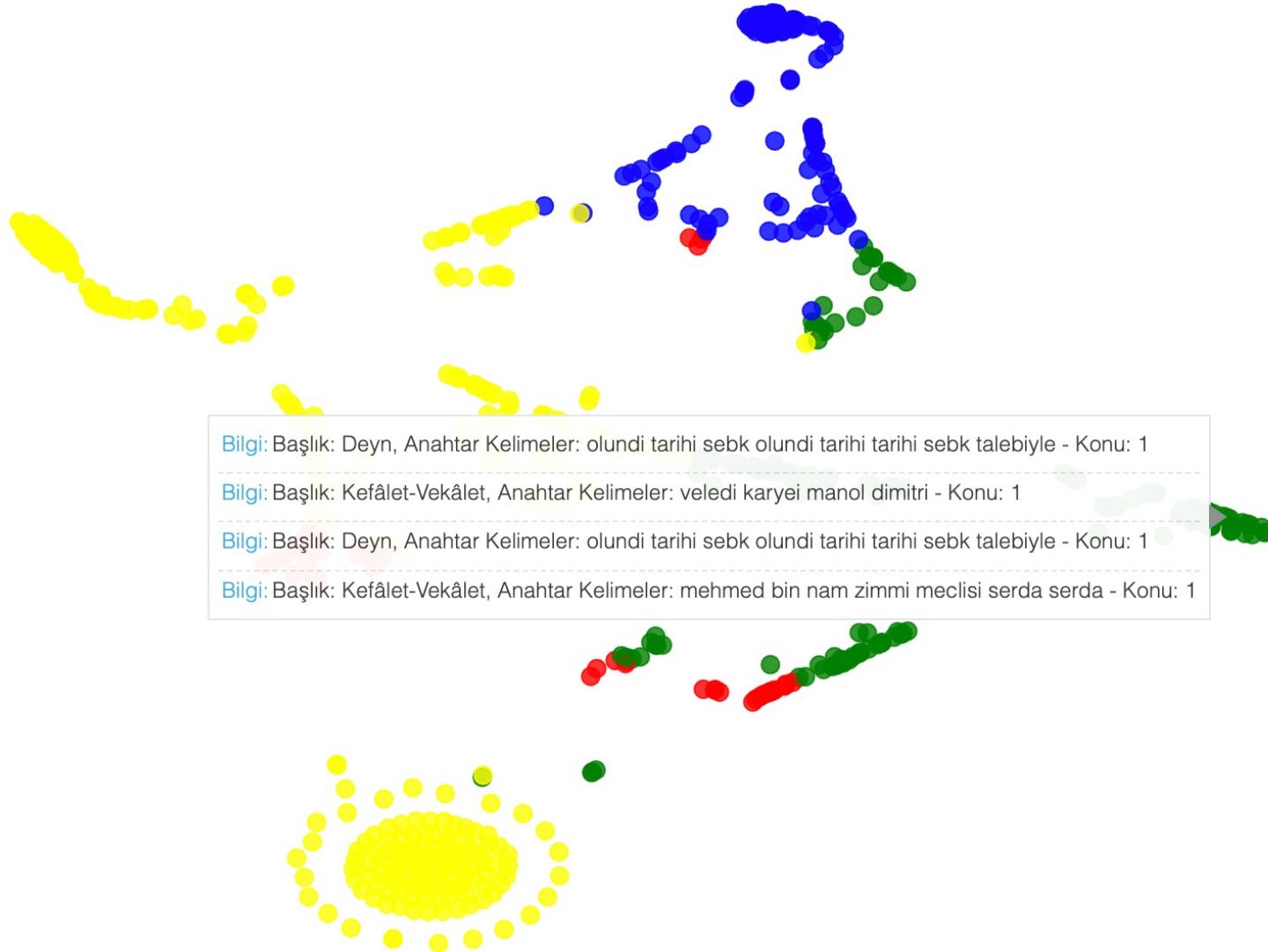


Konuların T-SNE görselleştirilmesi

- Topic-2
- Topic-4
- Topic-3
- Topic-1



- Topic-2
- Topic-4
- Topic-3
- Topic-1



Conclusion

- With this model, it has been seen more clearly which topics in the Kadi Registers are included with which words.
- This facilitates the identification of documents and automatic classification.
- In the transcription of documents, each case is assigned to a single subject.
- However, in the LDA model, it is possible to have different topics in each case at the same time.
- At the point of access to these documents, the user is offered a wider choice of topics.
- the LDA algorithm is based on the bag of words approach!! **(Limit)**

Perspectives on Data Access And Use at Scale: Lessons from the Field

Topic Modeling on the Kadi Registers

Dr. Sümeyye Akça

Assistant Professor
Marmara University
Information and Records Management
Turkey/Istanbul

sumeyyesakca@gmail.com

sumeyye.akca@marmara.edu.tr