

Abstract

Gensim is a pure Python library that fights on two fronts: 1) digital document indexing and similarity search; and 2) fast, memory-efficient, scalable algorithms for Singular Value Decomposition and Latent Dirichlet Allocation. The connection between the two is unsupervised, semantic analysis of plain text in digital collections. Gensim was created for large digital libraries, but its underlying algorithms for large-scale, distributed, online SVD and LDA are like the Swiss Army knife of data analysis—also useful on their own, outside of the domain of Natural Language Processing.

The Digital Library Stuff

Demo over arXiv: <http://aura.fi.muni.cz:8080/> (434,894 science articles).

DEVELOPED for the DML-CZ and EuDML¹ (digital math libraries) projects, as a vector-space alternative to the “find similar articles” functionality:

- Python server that runs as a daemon
 - Python/Java/C# clients (communication via Pyro/Pyrolite)
 - Clients **train a semantic model** on the server
 - Clients issue **add/remove/replace documents** requests
 - documents converted to “semantic” vectors using the model
 - Clients issue **queries for the most similar documents**
- AN eye on **performance** (numbers using my MacBookPro laptop C2D@2.53Ghz, vecLib for BLAS):
- Memory efficient **data streaming**
 - generators+iterators everywhere
 - train/index on corpora larger than RAM
 - Fast semantic model training (see to the right)
 - Efficient **incremental indexing**
 - 1.2k docs per minute, biggest part of it parsing and tokenizing input
 - At the lowest level, queries = matrix multiplications
 - index shards as NumPy&SciPy.sparse matrices mmap’ed from disk

The Math Stuff

STATISTICAL analysis of **co-occurrence patterns** to identify latent structure. In NLP: word co-occurrence over a corpus of plain text documents (no metadata).

- Training corpus as an implicit word-document matrix
 - sparse, much larger than RAM, streamed over sequentially
- Create a semantic model that captures corpus structure
- unique **Latent Semantic Analysis** (truncated SVD) and **Latent Dirichlet Allocation** implementations:
 - *one-pass*: each observation seen only once during training
 - *incremental*: can update model with new observations efficiently
 - *distributed*: can use Pyro to split the work over several machines/cores
 - *constant memory*: no $O(\#observations)$ required
 - ⇒ online training, can **process infinite data streams!**
- Using a trained model, can transform any plain text document to its “semantic” representation (see to the left)

EFFICIENCY: training LSA (truncated SVD) over the full version of English Wikipedia on my MBP laptop with vecLib BLAS:

- 3.5M docs, 100K vocab, 5.4G sparse non-zeroes
- Training: 400 factors in 6.5h
- Transforming: 18k docs/m using the 400-factor LSA model

Similar articles to article

CHEN, HUANYIN
Strong separativity over exchange rings. (English). Czechoslovak Mathematical Journal, vol. 58 (2008), issue 2, pp. 417-428

[-> Back to article](#)

Method LSI	Method RP	Method TFIDF
Generalized SVS -rings ...	Exchange rings with st...	Exchange rings with st...
ExchNUMDAM: Generalized SVS -rings and von Neumann regular rings	generalized SVS -rings ...	Exchange rings in whic...
Exchange rings in whic...	Exchange rings in whic...	Diagonal reductions of...
Rings which have proje...	Rings which have proje...	The least separative c...
Epimorphisms of regula...	S - ω $1S$ -generated u...	Note on the congruence...
Von Neumann regular ri...	Diagonal reductions of...	Extension of measure-l...
A general theory of Fo...	Von Neumann regular ri...	Integration in partial...
On SVS -rings and unit...	Steady ideals and rings	Modularity and distrib...
Extensions of S GM S -rings	Dualities over compact...	On abelian groups by w...
S ES-rings and differen...	The p. p. ring and the...	Extensions of S GM S -rings

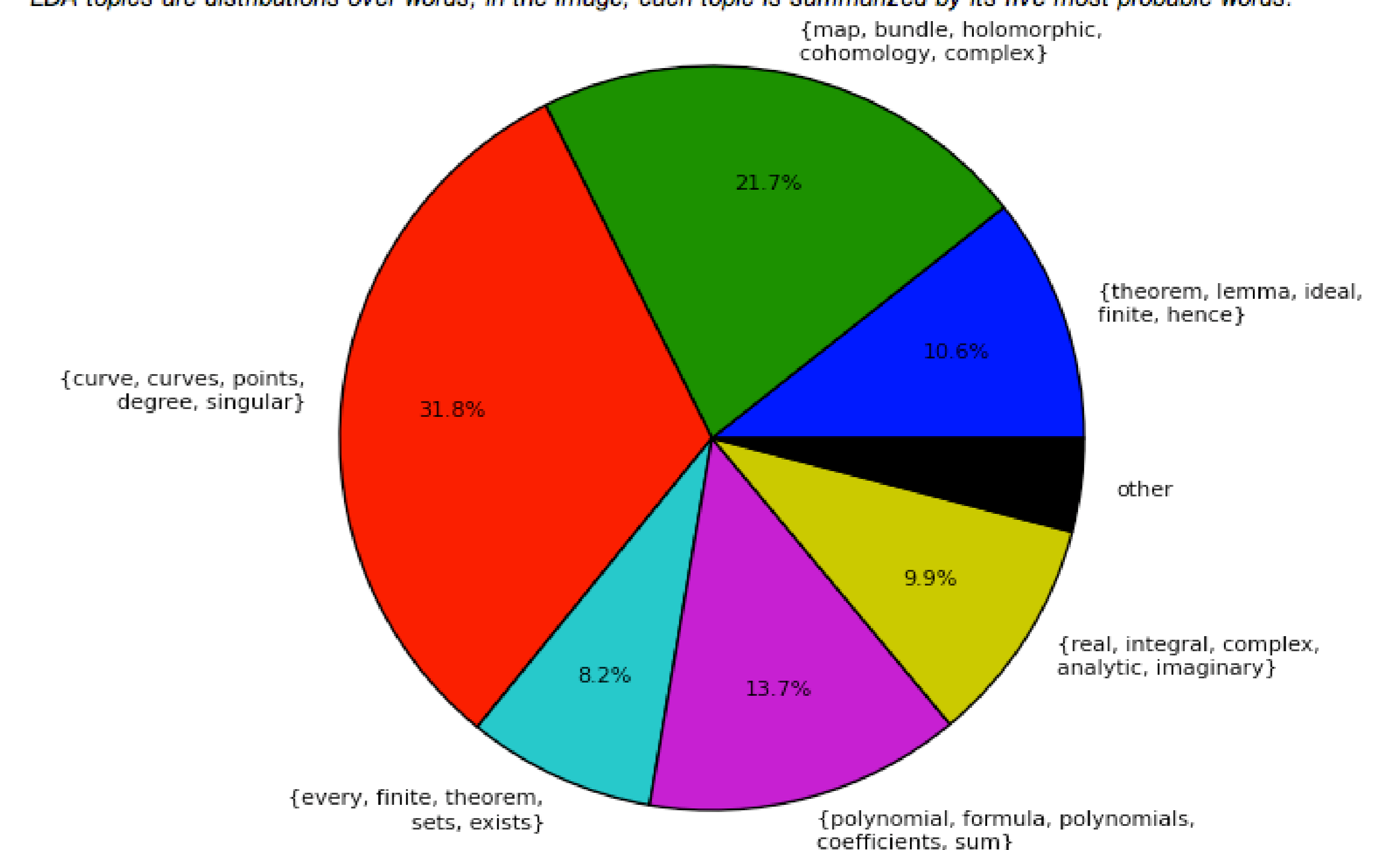
We appreciate your feedback to the methods which determine similarity of articles (e.g. which method is better, ...). Please [contact us](#). It will be helpful for future development.

[-> Back to article](#)

LDA Topics Pie Chart for math.0406240:

Each slice represents a different topic. The size of the slice corresponds to “how much is the article about this topic?”. Topics which contribute <6% to the above document are aggregated under “other”.

LDA topics are distributions over words; in the image, each topic is summarized by its five most probable words.



Credits

GENSIM is built on top of an excellent open-source Python stack: NumPy, SciPy and Pyro. Our work has been partially supported by the Ministry of Education of Czech Republic within the Center of Basic Research LC536 and by the European Union through its Competitiveness and Innovation Programme (Policy Support Programme, “Open access to scientific information”, Grant Agreement No. 250503). Many thanks to gensim contributors and testers. Gensim is licensed under LPGL—get it from PyPI or clone from github (just google it).

References

- [1] R. Řehůřek. Fast and Faster: A Comparison of Two Streamed Matrix Decomposition Algorithms. In *NIPS 2010 Workshop on Low-rank Methods for Large-scale Machine Learning*, Vancouver, Canada, 2010.
- [2] R. Řehůřek. Subspace Tracking for Latent Semantic Analysis. In *Advances in Information Retrieval*, volume 6611 of *Lecture Notes in Computer Science*, pages 289–300. Springer, 2011.
- [3] R. Řehůřek and P. Sojka. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta, May 2010.

¹The Czech Digital Mathematics Library <http://dm1.cz> and The European Digital Mathematics Library <http://www.eudml.eu>