

International Conference on Computational Social Science 2015

Topic Modeling Stability and Granulated LDA

Sergei Koltcov, Olessia Koltsova, Sergey Nikolenko.
<http://linis.hse.ru/en/>

WHY TOPIC MODELING STABILITY IS IMPORTANT FOR SOCIAL SCIENCE?

- Large text corpora (e.g. user generated content) are increasingly becoming an object of social science research.
- Social scientists assume that collections of texts are devoted to a number of topics in a certain proportion that reflect real interest of the text authors to those topics.
- **Social scientists also expect that mathematicians can provide them an instrument** that can detect these topics and their proportions.
- If an instrument provides “false” – e.g. biased or unstable – results, social scientists can not mine topics and draw any meaningful conclusions.

RESEARCH GOAL

In this work, we propose GLDA - a new, more stable modification of LDA, the most well-known topic modeling algorithm.

In this presentation:

- We explain the idea of GLDA based on word co-occurrence regularization.
- We report experimental results with three algorithms including GLDA.
- We explain a new metrics of stability - a modified version of Kullback – Leibler divergence that was used in the experiments along with traditional Jaccard coefficient.

Dataset

101,481 posts from the Russian LiveJournal.

172,939,000 tokens (unique words).

LDA Joint Probability

Terms: **D** – space of documents, **W** – space of unique words, **Z** – space of topics.
 Topics are hidden parameters, which have to be found in simulation. **Joint Probability** can be calculated according to the following equations:

$$p(w | z, \beta) = \int p(w | z, \Phi) p(\Phi | \beta) d\Phi \quad p(z | \alpha) = \int p(z | \Theta) p(\Theta | \alpha) d\Theta$$

Θ, Φ : words – topics and document – topics matrices. Those can be found out by different algorithms.



Mean field variational inference

$$L(\Phi, \Theta) = \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} \rightarrow \max$$



Collapsed Gibbs sampling

$$P(z_i = j | w_i = m, z_{-i}, w_{-i}) \approx \frac{C_{m,j}^{WT} + \beta}{\sum_m C_{m,j}^{WT} + V\beta} \cdot \frac{C_{d,j}^{DT} + \alpha}{C_{d,j}^{DT} + \alpha T}$$

In terms of matrix factorization LDA solution is:

$$F[\text{documents} \times \text{words}] = \Theta[\text{documents} \times \text{topics}] \cdot \Phi[\text{topics} \times \text{words}]$$

Mathematical vision of LDA

$$F[\text{documents} \times \text{words}] = \Theta[\text{documents} \times \text{topics}] \cdot \Phi[\text{topics} \times \text{words}]$$

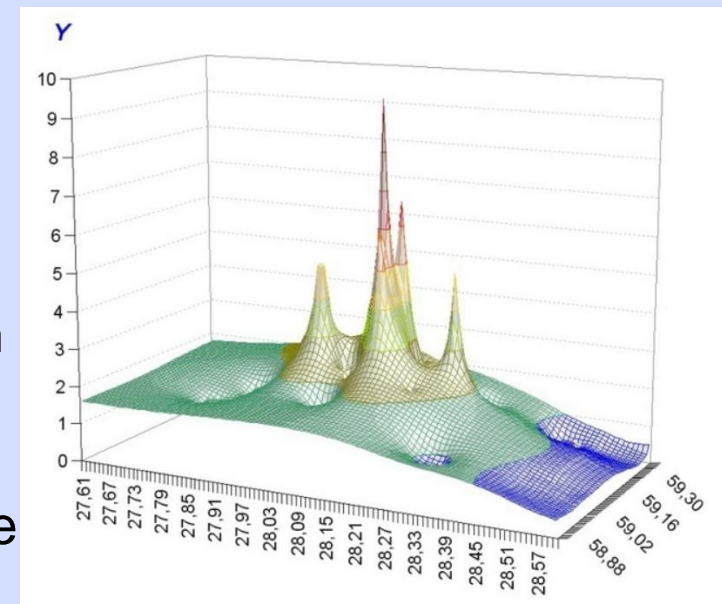
Matrix \mathbf{F} represents a dataset. Our dataset can be expressed in terms of two low dimension matrices. Process of sampling is the process of approximation of matrix \mathbf{F} by two matrices $\mathbf{\Phi}$ and $\mathbf{\Theta}$. **But:**

$$F = \Theta \cdot \Phi = (\Theta \cdot R) \cdot (R^{-1} \Phi) = \Theta' \cdot \Phi'$$

Matrix \mathbf{F} can be approximated by different combinations of matrices which means LDA has many solutions, e.i. many local minima..

Partly, the problem can be solved by using method of regularization whose main idea is in using additional information.

Regularization in topic modeling is a procedure that introduces reasonable restrictions on matrices.



Latent Dirichlet Allocation (Gibbs sampling)

$$P(z_i = j | w_i = m, z_{-i}, w_{-i}) \approx \frac{C_{m,j}^{WT} + \beta}{\sum_m C_{m,j}^{WT} + V\beta} \cdot \frac{C_{d,j}^{DT} + \alpha}{C_{d,j}^{DT} + \alpha T}$$

$C_{m,j}^{WT}$ - Matrix; cells: number of times a word \mathbf{w} was assigned to topic \mathbf{t} ,

$C_{d,j}^{DT}$ - Matrix; cells: number of times a word \mathbf{w} in document \mathbf{d} is assigned to topic \mathbf{t} .

$\sum_m C_{m,j}^{WT} = n_t$ - Vector; cells: number of words assigned to topic \mathbf{t} ,

$C_{d,j}^{DT} = n_d$ Length of document \mathbf{d} in words

Results of simulation:

1. Matrix of words distribution on topics. 2. Matrix of documents distribution on topics.

$$\theta_{dj} = \frac{C_{d,j}^{DT} + \alpha}{C_{d,j}^{DT} + T\alpha}$$

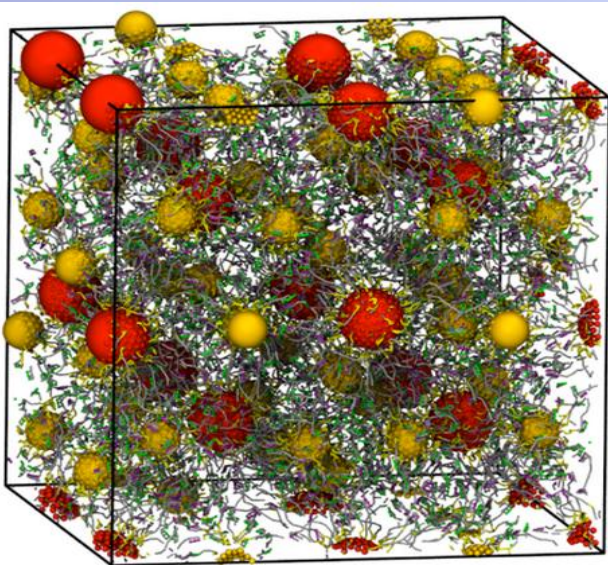
$$\phi_{m,j} = \frac{C_{m,j}^{WT} + \beta}{\sum_m C_{m,j}^{WT} + V\beta}$$

LDA is a pLSA model with coefficients of regularization (α, β).

Semi-Supervised Latent Dirichlet Allocation (Gibbs sampling)

Next level of regularization is based on the following idea. If we have initial distribution of words (anchor words) over topics, then we are able to fix or glue words to topics. Therefore, when the algorithm faces an anchor word during sampling, it does not change the connection between the topic and the word. But the other words are sampled according to the standard procedure.

$$p(z, w, \alpha, \beta) \propto \begin{cases} z = t & \xrightarrow{\text{Initial anchor words distribution}} \\ q(z, w, \alpha, \beta) & \xrightarrow{\text{Standard Gibbs sampling}} \end{cases}$$



The SLDA modeling behaves as a process of crystallization, where anchor words are centers of crystals. The words that often co-occur with anchor words stick together during simulation and form the body of topics.

Therefore the fixed lists of anchor words assigned to topics lead to stabilization of topic modeling.

But SLDA works good if the list of words is known beforehand which is not always the case in social science.

GRANULATED LDA (based on Gibbs sampling)

Granulated LDA based on idea that each document from a collection can be regarded as a granulated surface – a notion we borrow from physics. Here a granule is a set of words located near each other. Therefore we can arrange sampling by granules.

All words within one granule belong to one topic. Therefore scanning documents and assigning granules to topics, the algorithm favors the words located near each other.

DOCUMENT

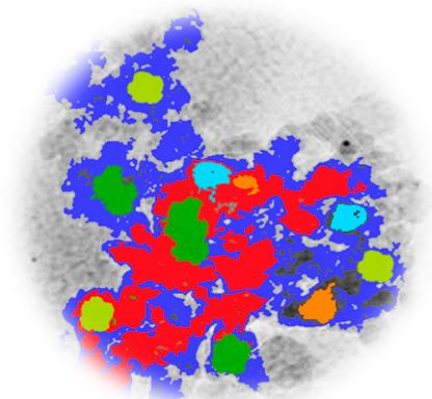
The central theme of **ethnic nationalists** is that «nations are defined by a shared heritage, which usually includes a **common language**, a common faith, and a **common ethnic ancestry**».[2] It also includes ideas of **a culture** shared between members of the group, and with their ancestors, and usually a shared language; however it is different from purely cultural definitions of «the nation» (which allow people to become members of a nation **by cultural assimilation**) and a purely linguistic definitions (which see «the nation» as all speakers of a specific language). Herodotus is the first who stated the main **characteristic of ethnicity**, with his famous account of what defines **Greek identity**, where he lists **kinship language, cults and customs**.

The central political tenet of **ethnic nationalism** is that **ethnic groups** can be identified unambiguously, and that each such group is entitled to **self-determination**.

The outcome of this right to **self-determination** may vary, from calls for self-regulated administrative bodies within an already-established society, to an **autonomous entity separate** from that **society**, to a **sovereign state** removed from **that society**. In international relations, it also leads to policies and movements for irredentism to claim a **common nation based upon ethnicity**.

KEYWORDS

ethnic	7
common language	3
culture	2
self-determination	2
society	3



granulated surface

GRANULATED LDA: Algorithm

Entrance: collection of documents D , number of topics $|T|$,
number of iterations, size of granules L ;

Initialization: $\varphi(w,t)$, $\theta(t,d)$ for all documents and topics $d \in D$, $w \in W$, $t \in T$;

Run external cycle along all documents (i)

Run internal cycle. Length of cycle is number words in documents i .

1. Generation of random number k . Max value of k is number of words in documents i .
2. Choosing word k from document i .
3. Calculating topic number t for word k .
4. Defining words that are around word k in document i .
5. Assigning topic t to all words which are within granule L .

End of internal cycle.

Updating the following matrices:

$$C_{m,j}^{WT} \quad C_{d,j}^{DT} \quad \sum_m C_{m,j}^{WT} = n_t$$

End of external cycle

Calculation of matrixes $\varphi(w,t)$, $\theta(t,d)$ based on

$$\theta_{dj} = \frac{C_{d,j}^{DT} + \alpha}{C_{d,j}^{DT} + T\alpha} \quad \phi_{m,j} = \frac{C_{m,j}^{WT} + \beta}{\sum_m C_{m,j}^{WT} + V\beta}$$

Visualization $\varphi(w,t)$, $\theta(t,d)$.

Evaluating LDA quality with Kullback–Leibler divergence and Jaccard coefficient

The **Kullback-Leibler divergence (K)** is a widely accepted distance measure between two probability distributions. Normalized **K** can be calculated according to the following formula.

$$Kn = \left(1 - \frac{K}{Max}\right) \cdot 100\% \quad \text{where} \quad Kn = 0.5 \sum_{k=1}^W \Omega_k^1 \log\left(\frac{\Omega_k^1}{\Omega_k^2}\right) + 0.5 \sum_{k=1}^W \Omega_k^2 \log\left(\frac{\Omega_k^2}{\Omega_k^1}\right)$$

IF $Kn=100\%$, then two topics are identical. IF $K=0$ then that topics are totally different.

Jaccard coefficient: $Jc=a/(a+b-k)$.

where **a** – one topic, **b** – another topic, **k** – number of words, which are identical in both topics.

$Jc = 1$, if two topic are identically, if **$Jc = 0$** then topic completely different.

Topic similarity thresholds

Level 90 - 93% (and more) means that first 50 words are almost identical.

Level about 85%: topics are completely different.

Similarity 0.935

USA	0.04734	USA	0.03567
American	0.02406	American	0.01804
Syria	0.02082	Syria	0.01758
Obama	0.01374	country	0.01495
weapon	0.01343	war	0.01361
war	0.01309	military	0.01246
president	0.01169	weapon	0.01084
UN	0.01018	Russia	0.01004
military	0.01014	Obama	0.00996
country	0.01005	president	0.0096
chemical	0.00944	UN	0.00869
Syrian	0.00851	international	0.00769

Similarity 0.854

USA	0.04734	water	0.01758
American	0.02406	help	0.01296
Syria	0.02082	city	0.01262
Obama	0.01374	far	0.01199
weapon	0.01343	house	0.01064
war	0.01309	east	0.0104
president	0.01169	region	0.00945
UN	0.01018	dam	0.0091
military	0.01014	flood	0.00904
country	0.01005	resident	0.00839
chemical	0.00944	injured	0.00714
Syrian	0.00851	FRS	0.00698

Words in topics are ordered according to probability.

Results of simulations

Dataset

1. 101481 posts from the Russian LiveJournal.
2. 172939 000 tokens (unique words).

Number of runs: **5 times**; number of topics: 200

Topic model (Gibbs sampling modification)	Number of stable topics According to Kullback- Leibler	Kullback–Leibler Divergence threshold (dictionary length= 1000 words)	Jaccard coefficient (100 most probable words in topic)
LDA	84	≈90%	≈0.37
SLDA	135	≈92.43%	≈0.16
GLDA (window ±1)	138	≈96%	≈0.6

CONCLUSION

1. We have proposed Granulated LDA as a regularized version of LDA.
2. For the purpose of detecting LDA stability, we have proposed a new topic similarity measure based on Kullback-Leibler divergence.
3. Based on this measure we have investigated stability of 3 modifications of LDA and have shown that Granulated LDA is more stable than other modifications.
4. Therefore, to be able to draw reliable sociological conclusions we recommend researchers either to run topic modeling several times, then extract stable topics that reappear across multiple runs or use Granulated LDA.

**The research is funded by the Basic Research Program
Of National Research University
Higher School of Economics**

THANK YOU!

← → ↻ 🏠 linis.hse.ru/en/o-nas



→ [HSE Campus in St. Petersburg](#) → [Internet Studies Lab](#) → About LINIS

RU EN

About LINIS

Laboratory for Internet Studies is a [team](#) of researchers from a variety of disciplines who have come together to do a pioneering research of issues that can not be addressed from within separate braches of science. Our mission is to study the Internet as a unique social, economic, linguistic and technical phenomenon. We also seek to learn more about the society that has “inhabited” the Internet and found it to be one of its natural modes of existence. For these goals we develop and adapt mathematical methods, algorithms and [software](#). All our [projects](#) are collective and include people with various skills and of different stages in their careers – from undergrads to PhD students to mature researchers. This enhances mutual enrichment and exchange of ideas. We support creative atmosphere and peer communication in our Lab; we welcome new initiatives and collaboration with business and policy makers and organize multiple [events](#). The results of our effort may be found in our [publications](#).