

Pólya Urn Latent Dirichlet Allocation

Alexander Terenin
Imperial College London

Joint work with Måns Magnusson,
Leif Jonsson, and David Draper

January 22, 2018

Talk for Carnegie Mellon University

arXiv:1704.03581

Latent Dirichlet Allocation

Latent Dirichlet Allocation

[DM Blei](#), [AY Ng](#), [MI Jordan](#) - [Journal of machine Learning research, 2003 - jmlr.org](#)

Abstract We describe latent Dirichlet allocation (LDA), a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modeled as a finite mixture over an underlying

☆  Cited by **21494** [Related articles](#) [All 127 versions](#)

The canonical topic model - everybody uses it!

Topic Modeling

Topic modeling is an area of natural language processing.

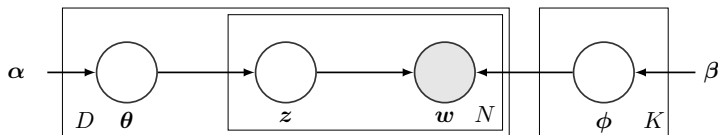
Setup

- We have a dataset consisting of *documents*.
- Each document is a collection of *word tokens*.

Goal: group documents into topics.

Approach: Bayesian learning.

Latent Dirichlet Allocation



Define a *topic* as a probability distribution over word tokens.

Assume each *document* has a probability distribution over topics.

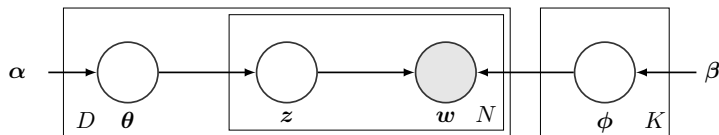
Assume documents are exchangeable, assume bag of words:

de Finetti's
theorem
 \implies

conditionally multinomial likelihood.

$\Phi_{K \times V}$: topic-word proportions. z : unobserved topic indicators.
 $\Theta_{D \times K}$: document-topic proportions. w : observed word tokens.

Latent Dirichlet Allocation



Φ : topic-word proportions. z : unobserved topic indicators.
 $K \times V$
 Θ : document-topic proportions. w : observed word tokens.
 $D \times K$

Need to learn Θ, Φ, z .

Bayesian Learning: (prior, likelihood) \rightarrow posterior.

$$\theta_k \sim \text{Dir}(\alpha) \qquad \phi_k \sim \text{Dir}(\beta)$$

Dirichlet distribution: conjugate with multinomial.

Training algorithm for Θ, Φ, z

We only care about z : integrate Θ, Φ out $\implies z \sim$ discrete.

Algorithm – sparse fully collapsed Gibbs sampling:

$$z_{i,d} \mid \mathbf{z}_{-i,d} \propto \frac{n_{k,v(i)}^{-i} + \beta}{n_{k,\cdot}^{-i} + V\beta} (m_{d,k}^{-i} + \alpha).$$

Many alternative methods: variational, geometric, spectral, etc.

MCMC approaches tend to give best-quality topics.

Problem: only uses sparsity in $\mathbf{m} \implies O\left[\sum_{i=1}^N K_{d_i}^{(\mathbf{m})}\right]$ per iteration.

Problem: completely sequential \implies slow.

Asynchronous Distributed LDA

Solution: pretend that

$$z_{i,d} \mid \mathbf{z}_{-i,d} \propto \frac{n_{k,v(i)}^{-i} + \beta}{n_{k,\cdot}^{-i} + V\beta} (m_{d,k}^{-i} + \alpha).$$

isn't sequential and just compute asynchronously.

Convergence analysis: much more complicated!¹

¹A. Terenin, D. Simpson, and D. Draper. Asynchronous Gibbs Sampling. In revision at Journal of the American Statistical Association: theory and methods. arXiv:1509.08999, 2017.

Partially Collapsed LDA

Can we do better? What if we don't integrate out Φ ?

Algorithm: sparse partially collapsed Gibbs sampling².

(1) Sample $\phi_k \sim \text{Dir}(\mathbf{n}_k + \beta)$.

(2) Sample $z_{i,d} \propto \phi_{k,v(i)}(m_{d,k}^{-i} + \alpha)$.

Step (1) is parallel by topics, step (2) is parallel by documents.

Empirical convergence rate is good: similar to fully collapsed LDA.

Problem: still only uses sparsity in $\mathbf{m} \implies O\left[\sum_{i=1}^N K_{d_i}^{(\mathbf{m})}\right]$.

Problem: requires $O(K \times V)$ memory.

²M. Magnusson, L. Jonsson, M. Villani, D. Broman. Sparse Partially Collapsed MCMC for Parallel Inference in Topic Models. Journal of Computational and Graphical Statistics, 2017.

Key idea

This work: completely eliminate both problems.

Observation: the problem is that $\boldsymbol{x} \sim \text{Dir}(\boldsymbol{\varpi})$ is a dense vector.

What if we made \boldsymbol{x} sparse?

We'd get a complexity of $O\left[\sum_{i=1}^N \min\left\{K_{d(i)}^{(\mathbf{m})}, K_{v(i)}^{(\Phi)}\right\}\right]$ – better!

Can we do that in a principled manner?

Dirichlet Processes and Pólya Urns

Pólya Urn: we have an urn with colored balls.

We take out one ball, and replace it with two balls of the same color.

Result: the de Finetti measure of a Pólya Urn is a Dirichlet Process.³

So, to sample from a Dirichlet process, simulate a large Pólya Urn.

Large in the sense of number of balls in the urn:

\implies concentration parameter of $DP(\varpi, F)$.

Henceforth: reparametrize $\text{Dir}(\varpi)$ as $\text{Dir}(\varpi, \mathbf{F})$ with $\varpi = \varpi \mathbf{F}$.

³D. Blackwell and J. MacQueen. Ferguson distributions via Pólya Urn schemes. The Annals of Statistics, 1973.

Pólya Urns and Bootstrap distributions

If we have 5 red balls and 2 green balls, what's the Pólya Urn's distribution? Exactly that: 5 red balls and 2 green balls:

$$\implies \text{MN}(\varpi, \mathbf{F}).$$

Problem: large multinomials are slow, cumbersome to work with.

Observation: same problem arises in distributed Bootstrap methods.

Idea: replace multinomial with Poisson process on an extended space.

The Poisson Pólya Urn distribution

Consider an n -dimensional probability vector x given by

$$x_i = \frac{\gamma_i}{\sum_{j=1}^n \gamma_j}$$

where $\gamma_i \sim \text{Pois}(\varpi_i)$. Call this the Poisson Pólya Urn distribution.

This is *almost* the Gamma representation of a Dirichlet distribution. We just replaced the Gammas with Poissons.

This is sparse and scalable. Does it actually work?

The Poisson Pólya Urn distribution

Theorem:⁴ if $\mathbf{x} \sim \text{PPU}(\varpi, \mathbf{F})$, and $\mathbf{x}^* \sim \text{Dir}(\varpi, \mathbf{F})$, we have

$$\|\mathbf{x} - \mathbf{x}^*\| \rightarrow 0$$

as $\varpi \rightarrow \infty$ in the Lévy-Prokhorov metric, and therefore both random vectors convergence in distribution.

In LDA, ϖ is driven by the number of word tokens, which is enormous.

Moreover, approximation is good even without asymptotics:

$$\mathbb{E}(x) = \mathbb{E}(x^*) \qquad \text{Var}(x) \cong \text{Var}(x^*).$$

⁴A. Terenin, M. Magnusson, L. Jonsson, and D. Draper. Pólya Urn Latent Dirichlet Allocation: a doubly sparse massively parallel sampler. In revision at IEEE Transactions on Pattern Analysis and Machine Intelligence. arXiv:1704.03581, 2017.

Pólya Urn LDA

Algorithm: Poisson Pólya Urn doubly sparse partially collapsed Gibbs sampling⁵.

(1) Sample $\phi_k \sim \text{PPU}(\mathbf{n}_k + \beta)$.

(2) Sample $z_{i,d} \propto \phi_{k,v(i)}(m_{d,k}^{-i} + \alpha)$.

Complexity: $O\left[\sum_{i=1}^N \min\left\{K_{d(i)}^{(\mathbf{m})}, K_{v(i)}^{(\Phi)}\right\}\right]$.

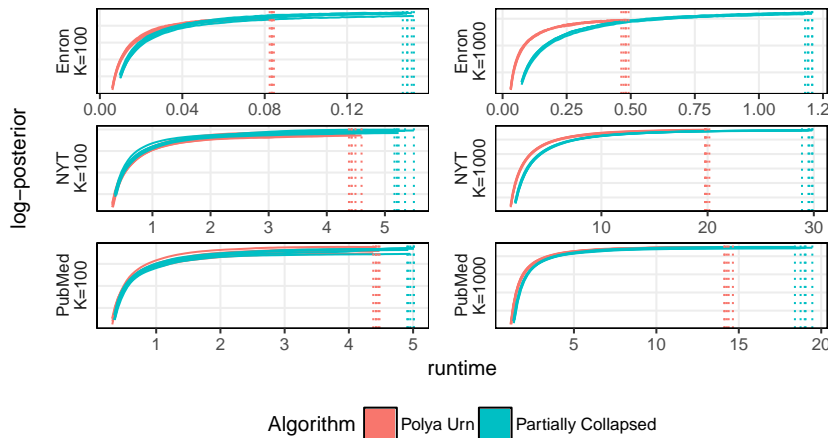
Everything is massively parallel.

Empirical convergence rate is good.

Same memory requirements as fully collapsed LDA.

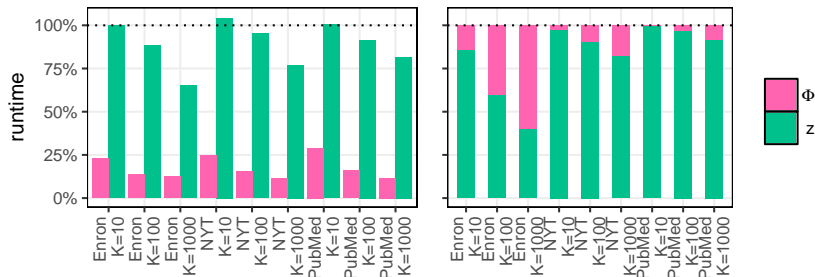
⁵A. Terenin, M. Magnusson, L. Jonsson, and D. Draper. Pólya Urn Latent Dirichlet Allocation: a doubly sparse massively parallel sampler. In revision at IEEE Transactions on Pattern Analysis and Machine Intelligence. arXiv:1704.03581, 2017.

Results



In a parallel setting, performance is state-of-the-art

Results

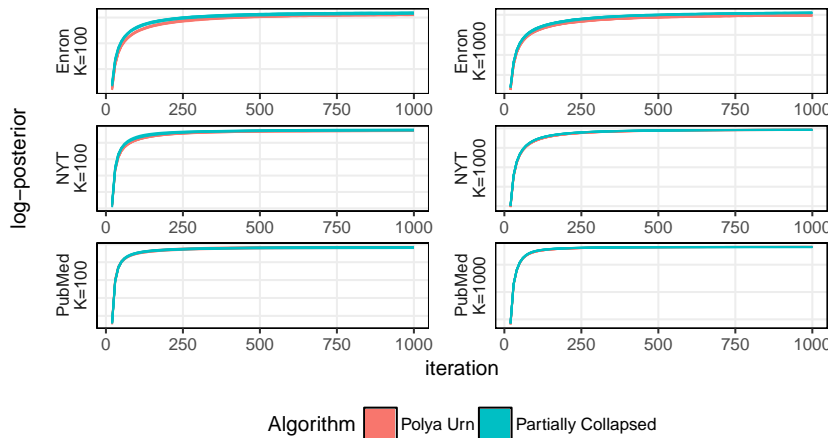


Where does the speedup come from?

Faster z due to additional sparsity.

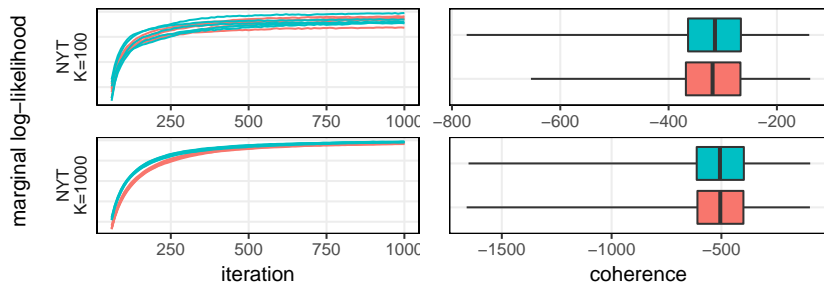
Faster Φ due to fast Poisson trick.

Results



Does it converge slower?

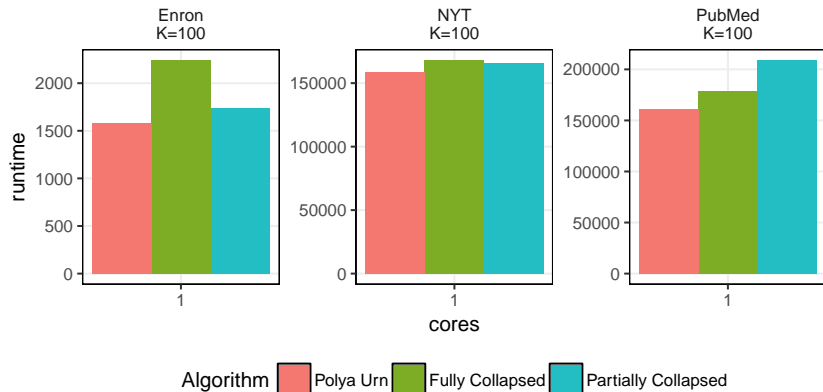
Results



Does it give good quality topics?

Yes: similar topic coherence and test set log marginal likelihood for z .

Results



Is single-core performance good?

Yes: also state-of-the-art among MCMC-based methods.

Results

Any downsides?

Approximation may not be good for small N and large K .

Sparse-on-sparse arithmetic is hard to implement well.

Conclusions

Poisson Pólya Urn: generic asymptotic approximation to Dirichlet.
Used in this work for doubly sparse massively parallel LDA sampler.

Theorem:⁶ if $\mathbf{x} \sim \text{PPU}(\varpi, \mathbf{F})$, and $\mathbf{x}^* \sim \text{Dir}(\varpi, \mathbf{F})$, we have

$$\|\mathbf{x} - \mathbf{x}^*\| \rightarrow 0$$

as $\varpi \rightarrow \infty$ in the Lévy-Prokhorov metric, and therefore both random vectors convergence in distribution.

Looks promising for other areas – try it yourself!

⁶A. Terenin, M. Magnusson, L. Jonsson, and D. Draper. Pólya Urn Latent Dirichlet Allocation: a doubly sparse massively parallel sampler. In revision at IEEE Transactions on Pattern Analysis and Machine Intelligence. arXiv:1704.03581, 2017.