

Rethinking LDA: moment matching for discrete ICA

Anastasia Podosinnikova, Francis Bach, and Simon Lacoste-Julien

INRIA - SIERRA project-team
École Normale Supérieure, Paris, France

Abstract

We consider moment matching techniques for estimation in Latent Dirichlet Allocation (LDA). By drawing explicit links between LDA and discrete versions of independent component analysis (ICA), we first derive a new set of cumulant-based tensors, with an improved sample complexity. Moreover, we reuse standard ICA techniques such as joint diagonalization of tensors to improve over existing methods based on the tensor power method. In an extensive set of experiments on both synthetic and real datasets, we show that our new combination of tensors and orthogonal joint diagonalization techniques outperforms existing moment matching methods.

1 Introduction

Topic models have emerged as flexible and important tools for the modelisation of text corpora. While early work has focused on graphical-model approximate inference techniques such as variational inference [6] or Gibbs sampling [21], tensor-based moment matching techniques have recently emerged as strong competitors due to their computational speed and theoretical guarantees [1, 3]. In this paper, we draw explicit links with the independent component analysis (ICA) literature (e.g. [19] and references therein) by showing a strong relationship between latent Dirichlet allocation (LDA) [6] and ICA [24, 25, 18]. We can then reuse standard ICA techniques and results, and derive new tensors with better sample complexity and new algorithms based on joint diagonalization.

2 Is LDA discrete PCA or discrete ICA?

Notation. Following the text modeling terminology, we define a corpus $X = \{x_1, \dots, x_N\}$ as a collection of N documents. Each document is a collection $\{w_{n1}, \dots, w_{nL_n}\}$ of L_n tokens. It is convenient to represent the ℓ -th token of the n -th document as a 1-of- M encoding with an indicator vector $w_{n\ell} \in \{0, 1\}^M$ with only one non-zero, where M is the vocabulary size, and each document as the count vector $x_n := \sum_{\ell} w_{n\ell} \in \mathbb{R}^M$. In such representation, the length L_n of the n -th document is $L_n = \sum_m x_{nm}$. We will always use index $k \in \{1, \dots, K\}$ to refer to topics, index $n \in \{1, \dots, N\}$ to refer to documents, index $m \in \{1, \dots, M\}$ to refer to words from the vocabulary, and index $\ell \in \{1, \dots, L_n\}$ to refer to tokens of the n -th document.

Latent Dirichlet allocation [6] is a generative probabilistic model for discrete data such as text corpora. In accordance to this model, the n -th document is modeled as an admixture over the vocabulary of M words with K latent topics. Specifically, the latent variable θ_n , which is sampled from the Dirichlet distribution, represents the topic mixture proportion over K topics for the n -th document. Given θ_n , the topic choice $z_{n\ell}|\theta_n$ for the ℓ -th token is sampled from the multinomial distribution with the probability vector θ_n . The token $w_{n\ell}|z_{n\ell}, \theta_n$ is then sampled from the multinomial distribution with the probability vector $d_{z_{n\ell}}$, or d_k if k is the index of the non-zero element in $z_{n\ell}$. This vector d_k is the k -th topic, that is a vector of probabilities over the words from the vocabulary subject to the simplex constraint, i.e. $d_k \in \Delta_M$, where $\Delta_M := \{d \in \mathbb{R}^M : d \succeq 0, \sum_m d_m = 1\}$. This generative process of a document (the index n is omitted for simplicity) can be summarized as

$$\begin{aligned}\theta &\sim \text{Dirichlet}(c), \\ z_{\ell}|\theta &\sim \text{Multinomial}(1, \theta), \\ w_{\ell}|z_{\ell}, \theta &\sim \text{Multinomial}(1, d_{z_{\ell}}),\end{aligned}\tag{1}$$

One can think of the latent variables z_ℓ as auxiliary variables which were introduced for convenience of inference, but can in fact be marginalized out, which leads to the following model

$$\begin{aligned}\theta &\sim \text{Dirichlet}(c), \\ x|\theta &\sim \text{Multinomial}(L, D\theta),\end{aligned}\tag{LDA model (2)}$$

where $D \in \mathbb{R}^{M \times K}$ is the topic matrix with the k -th column equal to the k -th topic d_k , and $c \in \mathbb{R}_{++}^K$ is the vector of parameters for the Dirichlet distribution. While a document is represented as a set of tokens w_ℓ in the formulation (1), the formulation (2) instead compactly represents a document as the count vector x . Although the two representations are equivalent, we focus on the second one in this paper and therefore refer to it as the LDA model.

Importantly, the LDA model does not model the length of documents. Indeed, although the original paper [6] proposes to model the document length as $L|\lambda \sim \text{Poisson}(\lambda)$, this is never used in practice and, in particular, the parameter λ is not learned. Therefore, in the way that the LDA model is typically used, it does not provide a complete generative process of a document as there is no rule to sample $L|\lambda$. In this paper, this fact is important, as we need to model the document length in order to make the link with discrete ICA.

Discrete PCA. Principal component analysis has the following probabilistic interpretation [29, 28]

$$\begin{aligned}\theta &\sim \text{Normal}(0, I_K), \\ x|\theta &\sim \text{Normal}(D\theta, \sigma^2 I_M),\end{aligned}\tag{3}$$

where $D \in \mathbb{R}^{M \times K}$ is a transformation matrix and σ is a parameter. As Buntine [9] mentions, the LDA model (2) can be seen as a discretization of the probabilistic PCA model (3) by replacing the normal likelihood with the multinomial one. By analogy with the self-conjugacy of the normal distribution, the Dirichlet prior is chosen as the conjugate prior for the multinomial distribution. Due to the close relation of models (2) and (3), LDA can be interpreted as a discrete PCA model.

Discrete ICA (DICA). Interestingly, a small extension of the LDA model allows its interpretation as a discrete independent component analysis (DICA) model. The extension naturally arises when the document length for the LDA model is modeled as a random variable from the gamma-Poisson mixture (which is equivalent to a negative binomial random variable), i.e. $L|\lambda \sim \text{Poisson}(\lambda)$ and $\lambda \sim \text{Gamma}(c_0, b)$, where $c_0 := \sum_k c_k$ is the shape parameter and $b > 0$ is the rate parameter. The LDA model (2) with such document length is equivalent (see Appendix A.1) to

$$\begin{aligned}\alpha_k &\sim \text{Gamma}(c_k, b), \\ x_m|\alpha &\sim \text{Poisson}([D\alpha]_m),\end{aligned}\tag{GP model (4)}$$

where all $\alpha_1, \alpha_2, \dots, \alpha_K$ are mutually independent, the parameters c_k coincide with the ones of the LDA model in (2), and the free parameter b can be seen (see Appendix A.2) as a scaling parameter for the document length when c_0 is already prescribed.

This model was originally introduced by Canny [11] and later named as a discrete ICA model [10]. It is more natural, however, to name model (4) as the gamma-Poisson (GP) model and the model

$$\begin{aligned}\alpha_1, \dots, \alpha_K &\sim \text{mutually independent}, \\ x_m|\alpha &\sim \text{Poisson}([D\alpha]_m)\end{aligned}\tag{DICA model (5)}$$

as the discrete ICA model. The only difference between (5) and the standard ICA model [24, 25, 18] (without additive noise) is the presence of the Poisson noise which enforces discrete, instead of continuous, values of x_m . Note also that (a) the discrete ICA model is a *semi-parametric* model that can adapt to any distribution on the topic intensities α_k and that (b) the GP model (4) is a particular case of both the LDA model (2) and the DICA model (5).

Thanks to this close connection between LDA and ICA, we can reuse standard ICA techniques to derive new efficient algorithms for topic modeling.

3 Moment matching for topic modeling

In general, the method of moments is based on the idea of estimating latent parameters of a probabilistic model by matching theoretical expressions of moments with their sample estimates. The

recent line of work by Anandkumar et al. [1, 3] discusses applications of the method of moments to several latent variable models including LDA, which results in computationally fast learning algorithms with theoretical guarantees. For LDA, their key ideas are (a) construction of the moments of the LDA model with some particular diagonal structure (called the “LDA moments” in the following) and (b) development of algorithms for estimating the model parameters by exploiting this particular diagonal structure. As discussed later, these algorithms are a particular kind of joint diagonalization algorithms on the sample estimates of expressions involving moments.

This paper has a similar high-level structure. In Section 3.1, we introduce novel cumulants for the GP/DICA models (called the “GP/DICA cumulants” in the following), which have a similar structure to the one of the LDA moments. This structure allows to reapply the algorithms of [1, 3] for the estimation of the model parameters, with the same theoretical guarantees. In addition, in Section 4, we consider another algorithm, which in turn is applicable to both the LDA moments and the GP/DICA cumulants. In Section 5, we experimentally compare these algorithms.

3.1 Cumulants of the GP and DICA models

In this section, we derive and analyze the novel cumulants of the DICA model. As the GP model is a particular case of the DICA model, all results of this section extend to the GP model.

The first three *cumulant tensors*¹ for the random vector x can be defined as follows

$$\text{cum}(x) := \mathbb{E}(x), \quad (6)$$

$$\text{cum}(x, x) := \text{cov}(x, x) = \mathbb{E}[(x - \mathbb{E}(x))(x - \mathbb{E}(x))^\top], \quad (7)$$

$$\text{cum}(x, x, x) := \mathbb{E}[(x - \mathbb{E}(x)) \otimes (x - \mathbb{E}(x)) \otimes (x - \mathbb{E}(x))], \quad (8)$$

where \otimes denotes the tensor product (see some properties of cumulants in Appendix B.1). The essential property of the cumulants (which does not hold for moments) that we use in this paper is that the cumulant tensor for a random vector with *independent* components is *diagonal*.

Let $y = D\alpha$; then for the Poisson random variable $x_m|y_m \sim \text{Poisson}(y_m)$, the expectation is $\mathbb{E}(x_m|y_m) = y_m$. Hence, by the law of total expectation and the linearity of expectation, the expectation in (6) has the following form

$$\mathbb{E}(x) = \mathbb{E}(\mathbb{E}(x|y)) = \mathbb{E}(y) = D\mathbb{E}(\alpha). \quad (9)$$

Further, the variance of the Poisson random variable x_m is $\text{var}(x_m|y_m) = y_m$ and, as x_1, x_2, \dots, x_M are conditionally independent given y , then their covariance matrix is diagonal, i.e. $\text{cov}(x, x|y) = \text{diag}(y)$. Therefore, by the law of total covariance, the covariance in (7) has the form

$$\begin{aligned} \text{cov}(x, x) &= \mathbb{E}[\text{cov}(x, x|y)] + \text{cov}[\mathbb{E}(x|y), \mathbb{E}(x|y)] \\ &= \text{diag}[\mathbb{E}(y)] + \text{cov}(y, y) = \text{diag}[\mathbb{E}(x)] + D\text{cov}(\alpha, \alpha)D^\top, \end{aligned} \quad (10)$$

where the last equality follows by the multilinearity property of cumulants (see Appendix B.1). Moving the first term from the RHS of (10) to the LHS, we define

$$S := \text{cov}(x, x) - \text{diag}[\mathbb{E}(x)]. \quad \text{DICA S-cum.} \quad (11)$$

From (10) and by the independence of $\alpha_1, \dots, \alpha_K$ (see Appendix B.3), S has the following diagonal structure

$$S = \sum_k \text{var}(\alpha_k) d_k d_k^\top = D \text{diag}[\text{var}(\alpha)] D^\top. \quad (12)$$

By analogy with the second order case, using the law of total cumulance, the multilinearity property of cumulants, and the independence of $\alpha_1, \dots, \alpha_K$, we derive in Appendix B.2 expression (22), similar to (10), for the third cumulant (8). Moving the terms in this expression, we define a tensor T with the following element

$$\begin{aligned} [T]_{m_1 m_2 m_3} &:= \text{cum}(x_{m_1}, x_{m_2}, x_{m_3}) + 2\delta(m_1, m_2, m_3)\mathbb{E}(x_{m_1}) \\ &\quad - \delta(m_2, m_3)\text{cov}(x_{m_1}, x_{m_2}) - \delta(m_1, m_3)\text{cov}(x_{m_1}, x_{m_2}) - \delta(m_1, m_2)\text{cov}(x_{m_1}, x_{m_3}), \end{aligned} \quad \text{DICA T-cum.} \quad (13)$$

¹The 2nd and 3rd cumulants coincide with the 2nd and 3rd central moments (but not at higher order).

where δ is the Kronecker delta. By analogy with (12) (Appendix B.3), the tensor T has the diagonal structure

$$T = \sum_k \text{cum}(\alpha_k, \alpha_k, \alpha_k) d_k \otimes d_k \otimes d_k. \quad (14)$$

In Appendix D.1, we recall (in our notation) the matrix S (37) and the tensor T (38) for the LDA model [1], which are analogues of the matrix S (11) and the tensor T (13) for the GP/DICA models. Slightly abusing terminology, we refer to the matrix S (37) and the tensor T (38) as the ‘‘LDA moments’’ and to the matrix S (11) and the tensor T (13) as the ‘‘GP/DICA cumulants’’. The diagonal structure (39) & (40) of the LDA moments is similar to the diagonal structure (12) & (14) of the GP/DICA cumulants, though arising through a slightly different argument, as discussed at the end of Appendix D.1. Importantly, due to this similarity, the algorithmic frameworks for both the GP/DICA cumulants and the LDA moments coincide.

The following sample complexity results apply to the sample estimates of the GP cumulants:²

Proposition 3.1. *Under the GP model, the expected error for the sample estimator \hat{S} (27) for the GP cumulant S (11) is:*

$$\mathbb{E} [\|\hat{S} - S\|_F] \leq \sqrt{\mathbb{E} [\|\hat{S} - S\|_F^2]} \leq O\left(\frac{1}{\sqrt{N}} \max[\Delta \bar{L}^2, \bar{c}_0 \bar{L}]\right), \quad (15)$$

where $\Delta := \max_k \|d_k\|_2^2$, $\bar{c}_0 := \min(1, c_0)$ and $\bar{L} := \mathbb{E}(L)$.

A high probability bound could be derived using concentration inequalities for Poisson random variables [7]; but the expectation already gives the right order of magnitude for the error (for example via Markov’s inequality). By following a similar analysis as in [2], we can rephrase the topic recovery error in term of the error on the GP cumulant. Importantly, the whitening transformation redivides the error on S (15) by \bar{L}^2 , which is the scale of S (see Appendix C.5 for details). This means that the contribution from \hat{S} to the recovery error will scale as $O(\frac{1}{\sqrt{N}} \max\{\Delta, \frac{\bar{c}_0}{\bar{L}}\})$, where both Δ and $\frac{\bar{c}_0}{\bar{L}}$ are smaller than 1 and can be very small. We do not present the exact expression for the expected squared error for the estimator of T , but due to a similar structure in the derivation, we expect the analogous bound of $\mathbb{E} [\|\hat{T} - T\|_F] \leq \frac{1}{\sqrt{N}} \max[\Delta^{3/2} \bar{L}^3, \bar{c}_0^{3/2} \bar{L}^{3/2}]$. In Appendix B.4, we present the expression (27) for an unbiased finite sample estimate \hat{S} of S and the expression (28) for an unbiased finite sample estimate \hat{T} of T . A sketch of a proof for Proposition 3.1 can be found in Appendix C.

Current sample complexity results of the LDA moments [1] can be summarized as $O(1/\sqrt{N})$. However, the proof (which can be found in the supplementary material [2]) analyzes only the case when finite sample estimates of the LDA moments are constructed from *one* triple per document, i.e. $w_1 \otimes w_2 \otimes w_3$ only, and not from the U-statistics that average multiple (dependent) triples per document as in the practical expressions (41) and (42). Moreover, one has to be careful when comparing upper bounds. Nevertheless, comparing the bound (15) with the current theoretical results for the LDA moments, we see that the GP/DICA cumulants sample complexity contains the ℓ_2 -norm of the columns of the topic matrix D in the numerator, as opposed to the $O(1)$ coefficient for the LDA moments. This norm can be significantly smaller than 1 for vectors in the simplex (e.g. $\Delta = O(1/\|d_k\|_0)$ for sparse topics). This suggests that the GP/DICA cumulants may have better finite sample convergence properties than the LDA moments and our experimental results in Section 5.2 are indeed consistent with this statement.

The GP/DICA cumulants have a somewhat more intuitive derivation than the LDA moments as they are expressed via the count vectors x (which are the sufficient statistics for the model) and not the tokens w_ℓ ’s. Note also that the construction of the LDA moments depend on the unknown parameter c_0 . Given that we are in an unsupervised setting and that moreover the evaluation of LDA is a difficult task [31], setting this parameter is non-trivial. In Appendix F.1, we investigate this dependence experimentally and observe that the LDA moments are somewhat sensitive to the choice of c_0 .

4 Diagonalization algorithms

How is the diagonal structure (12) of S (11) and (14) of T (13) going to be helpful for the estimation of the model parameters? This question has already been thoroughly investigated in the signal

²Note that the expected squared error for the DICA cumulants is similar, but the expressions are less compact and, in general, depend on the prior on α_k .

processing literature more than two decades ago (see, e.g., [12, 13, 15, 23, 16, 19] and references therein) and was recently brought back to the machine learning community (see [1, 3] and references therein), approach that we review in this section. Note that the algorithms of this section apply to both the LDA moments and the GP/DICA cumulants due to their similar diagonal structure.

For simplicity, let us rewrite expressions (12) and (14) for S and T as follows

$$S = \sum_k s_k d_k d_k^\top, \quad T = \sum_k t_k d_k \otimes d_k \otimes d_k, \quad (16)$$

where $s_k := \text{var}(\alpha_k)$ and $t_k := \text{cum}(\alpha_k, \alpha_k, \alpha_k)$. Introducing the rescaled topics $\tilde{d}_k := \sqrt{s_k} d_k$, we can also rewrite $S = \tilde{D} \tilde{D}^\top$. Following the same assumption from [1] that the topic vectors are linearly independent (and thus \tilde{D} has full rank), we can compute a whitening matrix $W \in \mathbb{R}^{K \times M}$ of S , i.e. a matrix such that $WSW^\top = I_K$ where I_K is the K -by- K identity matrix (see Appendix E.1 for more details). We then obtain that the vectors $z_k := W \tilde{d}_k$ form an orthonormal set of vectors.

Further, let us define a projection $\mathcal{T}(v) \in \mathbb{R}^{K \times K}$ of a tensor $\mathcal{T} \in \mathbb{R}^{K \times K \times K}$ onto a vector $u \in \mathbb{R}^K$:

$$\mathcal{T}(u)_{k_1 k_2} := \sum_{k_3} \mathcal{T}_{k_1 k_2 k_3} u_{k_3}. \quad (17)$$

Applying the multilinear transformation (see, e.g., [3] for the definition) with W^\top to the tensor T from (16) and then projecting the resulting tensor $\mathcal{T} := T(W^\top, W^\top, W^\top)$ onto some vector $u \in \mathbb{R}^K$, we obtain

$$\mathcal{T}(u) = \sum_k \tilde{t}_k \langle z_k, u \rangle z_k z_k^\top \quad (18)$$

where $\tilde{t}_k := t_k / s_k^{3/2}$ is due to the rescaling of topics and $\langle \cdot, \cdot \rangle$ stands for the inner product. As the vectors z_k are orthonormal, the pairs z_k and $\lambda_k := \tilde{t}_k \langle z_k, u \rangle$ can be seen as eigenpairs of the matrix $\mathcal{T}(u)$, which are uniquely defined if the eigenvalues λ_k are all different. If they are unique, we can recover the GP/DICA (as well as LDA) model parameters via $\tilde{d}_k = W^\dagger z_k$ and $\tilde{t}_k = \lambda_k / \langle z_k, u \rangle$.

This, in fact, is the spectral algorithm for LDA [1] and its predecessor, the fourth-order³ blind identification algorithm [12, 13]. Indeed, one can define finite sample estimates⁴ \hat{S} and \hat{T} of S (11) and T (13) and expect that they possess approximately the diagonal structure (12) and (14) and, therefore, the reasoning from above can be applied, under the assumption that the effect of the sampling error is controlled.

This spectral algorithm, however, is known to be quite unstable in practice (see, e.g., [14]). To overcome this problem, some other algorithms were proposed. The most notable ones are probably the FastICA algorithm [23] and the JADE algorithm [16]. The FastICA algorithm, with appropriate choice of a contrast function, estimates iteratively the topics, making use of the orthonormal structure (18), and performs the deflation procedure at every step. The recently introduced tensor power method (TPM) for the LDA model [3] is close to the FastICA algorithm. Alternatively, the JADE algorithm modifies the spectral algorithm by performing *multiple* projections for (18) and then jointly diagonalizing the resulting matrices with an orthogonal matrix. The spectral algorithm is a special case of this orthogonal joint diagonalization algorithm when only one projection is chosen. Importantly, a fast implementation [17] of the orthogonal joint diagonalization algorithm from [8] was proposed, which is based on closed-form iterative Jacobi updates (see, e.g., [27] for the later).

In practice, the orthogonal joint diagonalization (JD) algorithm is more robust than FastICA (see, e.g., [5, p.30]) or the spectral algorithm. Moreover, although the application of the JD algorithm for the learning of topic models was mentioned in the literature [3, 26], it was never implemented in practice. In this paper, we apply the JD algorithm for the diagonalization of the GP/DICA cumulants as well as the LDA moments, which is described in Algorithm 1. Note that the choice of a projection vector $v_p \in \mathbb{R}^M$ obtained as $v_p = \widehat{W}^\top u_p$ for some vector $u_p \in \mathbb{R}^K$ is important and corresponds to the multilinear transformation of \hat{T} with \widehat{W}^\top along the third mode. Importantly, in Algorithm 1, the joint diagonalization routine is performed over $(P+1)$ $K \times K$ matrices, where

³Note that (a) the factorization of $S = \tilde{D} \tilde{D}^\top$ does not uniquely determine \tilde{D} as one can equivalently use $S = (\tilde{D}U)(\tilde{D}U)^\top$ with any orthogonal $K \times K$ matrix U . Therefore, one has to consider higher than the second order information; (b) in ICA the fourth-order tensors are used, because the third cumulant of the Gaussian distribution is zero, which is not the case in the DICA/LDA models, where the third order information is sufficient.

⁴The precise expression for our finite sample estimates \hat{S} (27) and \hat{T} (28) of S (11) and T (13) for the GP/DICA cumulants is somewhat cumbersome and is therefore moved to Appendix E.2. For completeness, we also present the finite sample estimates \hat{S} (41) and \hat{T} (42) of S (37) and T (38) for the LDA moments (which are consistent with the ones suggested in [3]) in Appendix E.3.

the number of topics K is usually not too big. This makes the algorithm computationally fast (Appendix E discusses the computational complexity). The same is true for the spectral algorithm, but not for TPM.

Algorithm 1 Orthogonal joint diagonalization (JD) algorithm for GP/DICA cumulants (or LDA moments)

- 1: Input: $X \in \mathbb{R}^{M \times N}$, K , P (number of random projections); (and c_0 for LDA moments)
- 2: Compute sample estimate $\hat{S} \in \mathbb{R}^{M \times M}$ ((27) for GP/DICA / (41) for LDA in Appendix E)
- 3: Estimate whitening matrix $\hat{W} \in \mathbb{R}^{K \times M}$ of \hat{S} (see Appendix E.1)
- 4: option (a) Choose vectors $\{u_1, u_2, \dots, u_P\} \subseteq \mathbb{R}^K$ uniformly at random from the unit ℓ_2 -sphere and set $v_p = \hat{W}^\top u_p \in \mathbb{R}^M$ for all $p = 1, \dots, P$ ($P = 1$ yields spectral algorithm)
- 5: option (b) Choose vectors $\{u_1, u_2, \dots, u_P\} \subseteq \mathbb{R}^K$ as the canonical basis e_1, e_2, \dots, e_K of \mathbb{R}^K and set $v_p = \hat{W}^\top u_p \in \mathbb{R}^M$ for all $p = 1, \dots, K$
- 6: For $\forall p$, compute $B_p = \hat{W} \hat{T}(v_p) \hat{W}^\top \in \mathbb{R}^{K \times K}$ ((50) for GP/DICA / (52) for LDA; Appendix E)
- 7: Perform orthogonal joint diagonalization of matrices $\{\hat{W} \hat{S} \hat{W}^\top = I_K, B_p, p = 1, \dots, P\}$ (see [8] and [17]) to find an orthogonal matrix $V \in \mathbb{R}^{K \times K}$ and vectors $\{a_1, a_2, \dots, a_P\} \subset \mathbb{R}^K$ such that

$$V \hat{W} \hat{S} \hat{W}^\top V^\top = I_K, \text{ and } V B_p V^\top \approx \text{diag}(a_p), p = 1, \dots, P$$

- 8: Output: joint diagonalization matrix $A = V \hat{W}$ and values $a_p, p = 1, \dots, P$
-

In Section 5.1, we compare experimentally the performance of the spectral, JD, and TPM algorithms for the estimation of the parameters of the GP/DICA as well as LDA models. We are not aware of any experimental comparison of these algorithms in the LDA context. While already working on this manuscript, the JD algorithm was also independently analyzed by [26] in the context of tensor factorization for general latent variable models. However, [26] focused mostly on the comparison of approaches for tensor factorization and their stability properties, with brief experiments using a latent variable model related but not equivalent to LDA for community detection. In contrast, we provide a detailed experimental comparison in the context of LDA in this paper, as well as propose a novel cumulant-based estimator.

Model parameters recovery. Algorithm 1 outputs a joint diagonalization matrix $A \in \mathbb{R}^{K \times M}$ that has the property that AD should be approximately diagonal up to a permutation of the columns of D . The standard approach [1] of taking the pseudo-inverse of A to get an estimate of the topic matrix D has a problem that it does not preserve the simplex constraint of the topics (in particular, the non-negativity of \tilde{D}). Due to the space constraints, we do not discuss this issue here, but we observed experimentally that this can potentially significantly deteriorate performance of all moment matching algorithms for LDA. We made an attempt to solve this problem by integrating the non-negativity constraint into the Jacobi-updates procedure of the orthogonal joint diagonalization algorithm, but the obtained results did not lead to any significant improvement. Therefore, for our experiments, we estimate the topic matrix by thresholding the negative values of the pseudo-inverse of A : $\hat{d}_k := \max(0, [A^\dagger]_{:,k}) / \max(0, [A^\dagger]_{:,k})_1$, where $[A^\dagger]_{:,k}$ is the k -th column of the pseudo-inverse A^\dagger of A (see Appendix E.4 for more details on the recovery of the model parameters for the GP model), and leave this issue as an open question for future research.

5 Experiments

In this section, (a) we compare experimentally the GP/DICA cumulants with the LDA moments and (b) we compare experimentally the spectral algorithm [1], the tensor power method [3] (TPM), the orthogonal joint diagonalization (JD) algorithm from Algorithm 1, and the variational inference algorithm for LDA [6].

Real data: the associated press (AP) dataset, from D. Blei’s web page⁵, with $N = 2,243$ documents and $M = 10,473$ words in the vocabulary and the average document length $\hat{L} = 194$; the NIPS papers dataset⁶ [20] of $N = 2,483$ NIPS papers and $M = 14,036$ words in the vocabulary, and the average document length $\hat{L} = 1,321$; the KOS blog entries dataset⁷, from the UCI Repository, with

⁵<http://www.cs.columbia.edu/~blei/lda-c/index.html>

⁶<http://ai.stanford.edu/~gal/data.html>

⁷<https://archive.ics.uci.edu/ml/datasets/Bag+of+Words>

$N = 3,430$ documents and $M = 6,906$ words in the vocabulary, and the average document length $\hat{L} = 136$. As the LDA moments require at least three tokens in each document, 1 document from the NIPS dataset and 3 documents from the AP dataset, which did not fulfill this requirement, were removed.

Semi-synthetic data are constructed by analogy with [4] and provide ground truth information for evaluation. First, the LDA parameters D and c are learned from the real datasets with the variational inference LDA and, then, toy data are sampled from a model of interest with the given parameters D and c . For each setting, data are sampled 5 times and the results are averaged. We plot error bars that are the minimum and maximum values. This provides the ground truth parameters D and c . For the AP data, $K \in \{10, 50\}$ topics are learned and, for the NIPS data, $K \in \{10, 90\}$ topics are learned. For larger K , the obtained topic matrix is ill-conditioned, which violates the identifiability condition for topic recovery using moment matching techniques [1]. All the documents with less than 3 tokens were resampled.

Sampling techniques. All the sampling models have the parameter c which is set to $c = c_0 \bar{c} / \|\bar{c}\|_1$, where \bar{c} is the learned c from the real dataset with variational LDA, and c_0 is a parameter that we can vary. The GP data are sampled from the gamma-Poisson model (4) with $b = c_0 / \hat{L}$ so that the expected document length is \hat{L} (see Appendix A.2). The LDA-fix(L) data are sampled from the LDA model (2) with the document length being fixed to a given L . The LDA-fix2(γ, L_1, L_2) data are sampled as follows: $(1 - \gamma)$ -portion of the documents are sampled from the LDA-fix(L_1) model with a given document length L_1 and γ -portion of the documents are sampled from the LDA-fix(L_2) model with a given document length L_2 .

Evaluation. Evaluation of topic recovery for semi-synthetic data is performed with the ℓ_1 -error between the recovered \hat{D} and true D topic matrices with the best permutation of columns:

$$\text{err}_{\ell_1}(\hat{D}, D) := \min_{\pi \in \text{PERM}} \frac{1}{2K} \sum_k \|\hat{d}_{\pi_k} - d_k\|_1 \in [0, 1]. \quad (19)$$

The minimization is over the possible permutations $\pi \in \text{PERM}$ of the columns of \hat{D} and can be efficiently obtained with the Hungarian algorithm for bipartite matching. For the evaluation of topic recovery in the real data case, we use an approximation of the log-likelihood for held out documents as the metric. The approximation is computed using a Chib-style method as described by [31] using the authors' implementation⁸. Note that this evaluation method is also applicable for the GP model, as it is a particular case of the LDA model.

Code and complexity. We used our own Matlab implementations of the GP/DICA cumulants, the LDA moments, the spectral algorithm, and the tensor power method, as, to our knowledge, no efficient implementation of these algorithms was available for LDA. The expressions (50) and (52) provide efficient formulas for fast computation of the GP/DICA cumulants and LDA moments ($O(RNK)$, where R is the largest number of non-zeros in the count vector x over all documents), which makes even the Matlab implementation fast for large datasets. For the orthogonal joint diagonalization algorithm, we implemented a faster C++ version of the previous Matlab implementation by J.-F. Cardoso. For variational inference, we used the code of D. Blei and modified it for the estimation of a non-symmetric Dirichlet prior c , which is known to be important [30]. The default values of the tolerance/maximum number of iterations parameters are used for variational inference. The computational complexity of one iteration for one document of the variational inference algorithm is $O(RK)$, where R is the number of non-zeros in the count vector for this document, which is then performed a significant number of times. Each experiment was run in a single thread.

Note that (a) for the large vocabulary size M , the computation of a whitening matrix can be expensive (in terms of both memory and time) and (b) the bottle-neck for the spectral, JD, and TPM algorithms is the computation of the cumulants/moments. One possible solution for (a) is to reduce the vocabulary size with, e.g., TF-IDF score, which is a standard practice in the topic modeling context. Another option is using a stochastic eigendecomposition (see, e.g., [22]) to approximate the whitening matrix. For (b): the spectral algorithm estimates the cumulants/moments only once and, therefore, is fast; joint estimation of P cumulants/moments for JD can be (and is) implemented much faster than estimation of P cumulants/moments by precomputing and reusing terms (e.g. WX) which appear in all cumulants/moments; for TPM, some parts of the cumulants/moments can also be precomputed, but as TPM normally does many more iterations than P , it is significantly slower. Note that the number of random restarts for TPM within one deflation step is set to 10 and

⁸<http://homepages.inf.ed.ac.uk/imurray2/pub/09etm/>

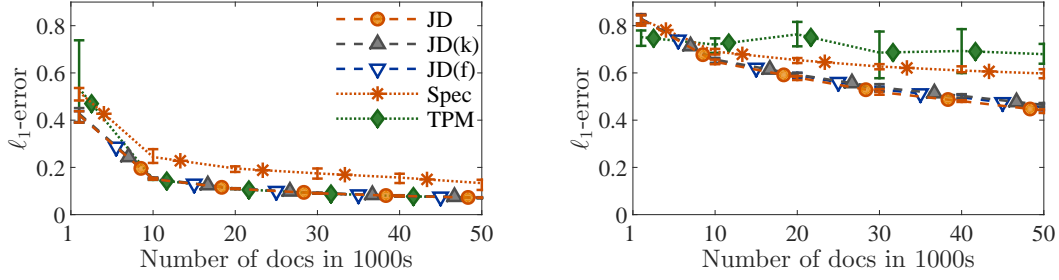


Figure 1: Comparison of the diagonalization algorithms. The topic matrix D and Dirichlet parameter c are learned for $K = 50$ from the AP dataset; c is scaled to sum up to 0.5 and b is set to fit the expected document length $\hat{L} = 200$. The semi-synthetic dataset is sampled from the GP model; number of documents N varies from 1,000 to 50,000. **Left:** GP/DICA moments. **Right:** LDA moments. *Note:* a smaller value of the ℓ_1 -error is better.

the maximum number of iterations for every run is set to 100; the run with the best objective is chosen. It is known that the runs which converge to a good solution converge fast [3].

Parameter c_0 for LDA. The construction of the LDA moments requires the parameter c_0 . For the semi-synthetic experiments, the true value of c_0 is provided to the algorithms. It means that the LDA moments, in this case, have access to some oracle information, which in practice is never available. For real data experiments, c_0 is set to the value obtained with variational inference. Experiments in Appendix F.1 show that this choice was somewhat important, however, this requires more thorough investigation.

5.1 Comparison of the diagonalization algorithms

In Figure 1, we give a comparison of the diagonalization algorithms on the semi-synthetic AP dataset for $K = 50$ using the GP model for sampling. We compare the tensor power method (TPM) [3], the spectral algorithm (Spec), the orthogonal joint diagonalization algorithm (JD) described in Algorithm 1 with different options to choose the random projections: JD(k) takes $P = K$ vectors u_p sampled uniformly from the unit ℓ_2 -sphere in \mathbb{R}^K and selects $v_p = W^\top u_p$ (option (a) in Algorithm 1); JD selects the full basis e_1, \dots, e_K in \mathbb{R}^K and sets $v_p = W^\top e_p$ (as JADE [16]) (option (b) in Algorithm 1); $JD(f)$ chooses the full canonical basis of \mathbb{R}^M as the vectors to project onto (is computationally expensive).

Although both the GP/DICA cumulants and LDA moments are well-specified for sampling from the GP model, the LDA moments have a slower finite sample convergence and, hence, a larger estimation error for the same value N . As expected, the spectral algorithm is always slightly inferior to the joint diagonalization algorithms. With the GP/DICA cumulants, where the estimation error is low, all algorithms demonstrate good performance, which also fulfills our expectations. However, although TPM shows almost perfect performance in the case of the GP/DICA cumulants (left), it significantly deteriorates for the LDA moments (right), which can be explained by the larger estimation error of the LDA moments and lack of robustness of TPM. For $N = 50,000$, the algorithms have the following runtimes (min, mean, max) in sec: JD-GP (140, 197, 267), JD(k)-GP (121, 176, 269), JD(f)-GP (1935, 2114, 2259), Spec-GP (87, 106, 125), TPM-GP (1857, 1993, 2102), JD-LDA (238, 308, 447), JD(k)-LDA (236, 292, 441), JD(f)-LDA (2998, 3547, 4939), Spec-LDA (89, 106, 147), TPM-LDA (1490, 2116, 2796). Note, due to random restarts, two runs of TPM/Spec are never the same. Computation of a whitening matrix is roughly 30 sec (this time is the same for all algorithms and is included in the numbers above). Overall, the orthogonal joint diagonalization algorithm with initialization of random projections as W^\top multiplied with the canonical basis in \mathbb{R}^K (JD) is both computationally efficient and fast.

5.2 Comparison of the GP/DICA cumulants and the LDA moments

In Figure 2, when sampling from the GP model (top, left), both the GP/DICA cumulants and LDA moments are well specified, which implies that the approximation error is low for both. The GP/DICA cumulants achieve low values of the estimation error already for $N = 10,000$ documents

independently of the number of topics, while the convergence is slower for the LDA moments. When sampling from the *LDA-fix(200)* model (top, right), the GP/DICA cumulants are mis-specified and their approximation error is high, although the estimation error is low due to the faster finite sample convergence. One reason of poor performance of the GP/DICA cumulants, in this case, is the absence of variance in document length. Indeed, if documents with two different lengths are mixed by sampling from the *LDA-fix2(0.5,20,200)* model (bottom, left), the GP/DICA cumulants' performance improves. Moreover, the experiment with a changing fraction γ of documents (bottom, right) shows that a non-zero variance on the length improves the performance of the GP/DICA cumulants. As in practice real corpora usually have a non-zero variance for the document length, this bad scenario for the GP/DICA cumulants is not likely to happen.

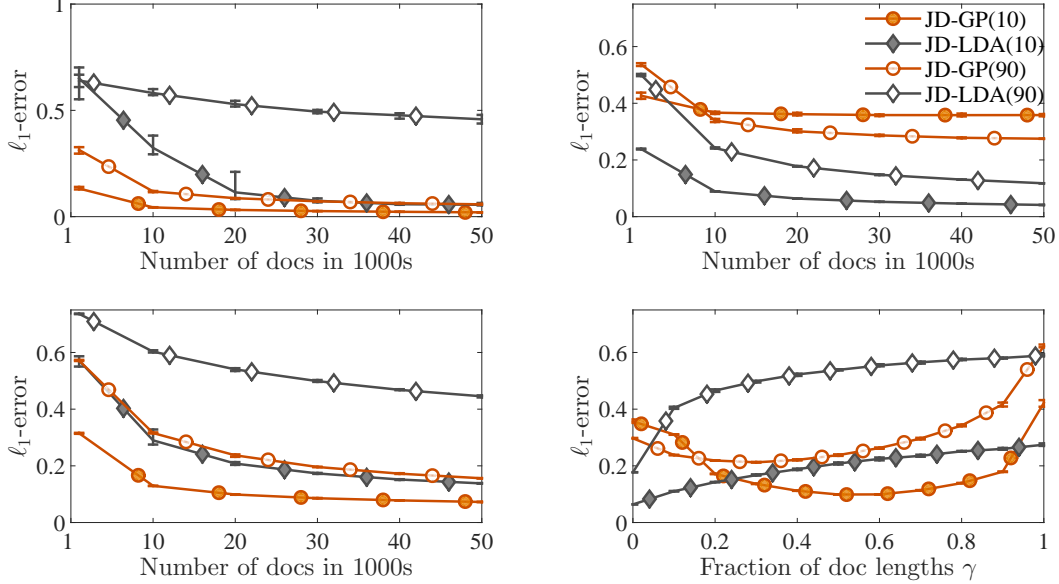


Figure 2: Comparison of the GP/DICA cumulants and LDA moments. Two topic matrices and parameters c_1 and c_2 are learned from the NIPS dataset for $K = 10$ and 90 ; c_1 and c_2 are scaled to sum up to $c_0 = 1$. Four corpora of different sizes N from 1,000 to 50,000: **top, left:** b is set to fit the expected document length $\hat{L} = 1300$; sampling from the *GP* model; **top, right:** sampling from the *LDA-fix(200)* model; **bottom, left:** sampling from the *LDA-fix2(0.5,20,200)* model. **Bottom, right:** the number of documents here is fixed to $N = 20,000$; sampling from the *LDA-fix2(γ ,20,200)* model varying the values of the fraction γ from 0 to 1 with the step 0.1. Note: a smaller value of the ℓ_1 -error is better.

5.3 Real data experiments

Each dataset is separated into 5 training/evaluation pairs, where the documents for evaluation are chosen randomly and non-repetitively among the folds (400 documents are held out for AP; 600 documents are held out for KOS). Then, the model parameters are learned for a different number of topics. The evaluation of the held-out documents is performed with averaging over 5 folds. In Figure 3, on the y-axis, the predictive log-likelihood in bits averaged per token is presented. JD-GP, Spec-GP, JD-LDA, and Spec-LDA are compared with variational inference (VI) and with variational inference initialized with the output of JD-GP (VI-JD). The orthogonal joint diagonalization algorithm with the GP/DICA cumulants (JD-GP) demonstrates competitive performance. In particular, the GP/DICA cumulants significantly outperform the LDA moments, and are better than variational inference. Interestingly, using variational inference after the topics have been learned by moment matching decreased performance in some cases. In Appendix F.3, a similar experiment for the NIPS dataset is presented.

6 Conclusion

In this paper, we have proposed a new set of tensors for a discrete ICA model related to LDA, where word counts are directly modelled. These moments make fewer assumptions regarding distributions, and are theoretically and empirically more robust than previously proposed tensors for LDA, both on synthetic and real data. Following the ICA literature, we showed that our joint diagonalization

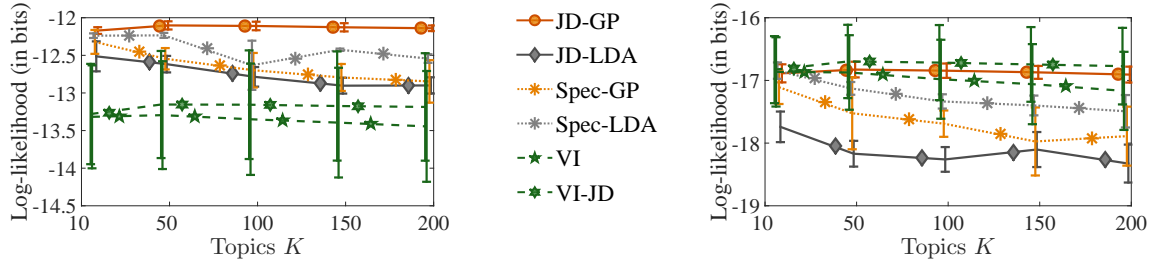


Figure 3: Experiments with real data. **Left:** the AP dataset. **Right:** the KOS dataset. *Note:* a higher value of the log-likelihood is better.

procedure is also more robust. Once the topic matrix has been estimated, it would be interesting to learn the unknown distributions of the independent topic intensities.

Acknowledgements. This work was partially supported by the MSR-Inria Joint Center.

References

- [1] A. Anandkumar, D.P. Foster, D. Hsu, S.M. Kakade, and Y.-K. Liu. A spectral algorithm for latent Dirichlet allocation. In *NIPS*, 2012.
- [2] A. Anandkumar, D.P. Foster, D. Hsu, S.M. Kakade, and Y.-K. Liu. A spectral algorithm for latent Dirichlet allocation. *CoRR*, abs:1204.6703, 2013.
- [3] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *J. Mach. Learn. Res.*, 15:2773–2832, 2014.
- [4] S. Arora, R. Ge, Y. Halpern, D. Mimno, A. Moitra, D. Sontag, Y. Wu, and M. Zhu. A practical algorithm for topic modeling with provable guarantees. In *ICML*, 2013.
- [5] F.R. Bach and M.I. Jordan. Kernel independent component analysis. *J. Mach. Learn. Res.*, 3:1–48, 2002.
- [6] D.M. Blei, A.Y. Ng, and M.I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:903–1022, 2003.
- [7] S. Boucheron, G. Lugosi, and P. Massart. *Concentration inequalities: a nonasymptotic theory of independence*. Oxford University Press, 2013.
- [8] A. Bunse-Gerstner, R. Byers, and V. Mehrmann. Numerical methods for simultaneous diagonalization. *SIAM J. Matrix Anal. Appl.*, 14(4):927–949, 1993.
- [9] W.L. Buntine. Variational extensions to EM and multinomial PCA. In *ECML*, 2002.
- [10] W.L. Buntine and A. Jakulin. Applying discrete PCA in data analysis. In *UAI*, 2004.
- [11] J. Canny. GaP: a factor model for discrete data. In *SIGIR*, 2004.
- [12] J.-F. Cardoso. Source separation using higher order moments. In *ICASSP*, 1989.
- [13] J.-F. Cardoso. Eigen-structure of the fourth-order cumulant tensor with application to the blind source separation problem. In *ICASSP*, 1990.
- [14] J.-F. Cardoso. High-order contrasts for independent component analysis. *Neural Comput.*, 11:157–192, 1999.
- [15] J.-F. Cardoso and P. Comon. Independent component analysis, a survey of some algebraic methods. In *ISCAS*, 1996.
- [16] J.-F. Cardoso and A. Souloumiac. Blind beamforming for non Gaussian signals. In *IEEE Proceedings-F*, 1993.
- [17] J.-F. Cardoso and A. Souloumiac. Jacobi angles for simultaneous diagonalization. *SIAM J. Mat. Anal. Appl.*, 17(1):161–164, 1996.
- [18] P. Comon. Independent component analysis, a new concept? *Signal Process.*, 36:287–314, 1994.
- [19] P. Comon and C. Jutten, editors. *Handbook of blind source separation: independent component analysis and applications*. Academic Press, 2010.

- [20] A. Globerson, G. Chechik, F. Pereira, and N. Tishby. Euclidean embedding of co-occurrence data. *J. Mach. Learn. Res.*, 8:2265–2295, 2007.
- [21] T. Griffiths. Gibbs sampling in the generative model of latent Dirichlet allocation. Technical report, Stanford University, 2002.
- [22] N. Halko, P.-G. Martinsson, and J.A. Tropp. Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):217–288, 2011.
- [23] A. Hyvärinen. Fast and robust fixed-point algorithms for independent component analysis. *IEEE Trans. Neural Netw.*, 10(3):626–634, 1999.
- [24] C. Jutten. *Calcul neuromimétique et traitement du signal: analyse en composantes indépendantes*. PhD thesis, INP-USM Grenoble, 1987.
- [25] C. Jutten and J. Héroult. Blind separation of sources, part I: an adaptive algorithm based on neuromimetic architecture. *Signal Process.*, 24:1–10, 1991.
- [26] V. Kuleshov, A.T. Chaganty, and P. Liang. Tensor factorization via matrix factorization. In *AISTATS*, 2015.
- [27] J. Nocedal and S.J. Wright. *Numerical optimization*. Springer, 2nd edition, 2006.
- [28] S. Roweis. EM algorithms for PCA and SPCA. In *NIPS*, 1998.
- [29] M.E. Tipping and C.M. Bishop. Probabilistic principal component analysis. *J. R. Stat. Soc.*, 61:611–622, 1999.
- [30] H.M. Wallach, D. Mimno, and A. McCallum. Rethinking LDA: why priors matter. In *NIPS*, 2009.
- [31] H.M. Wallach, I. Murray, R. Salakhutdinov, and D. Mimno. Evaluation methods for topic models. In *ICML*, 2009.

A Appendix. The GP model

A.1 The connection between the LDA and GP models

To show that the LDA model (2) with the additional assumption that the document length is modeled as a gamma-Poisson random variable is equivalent to the GP model (4), we show that:

- when modeling the document length L as a Poisson random variable with parameter λ , the count vectors x_1, x_2, \dots, x_M are mutually independent Poisson random variables;
- the Gamma prior on λ reveals the connection $\alpha_k = \lambda\theta_k$ between the Dirichlet random variable θ and the mutually independent gamma random variables $\alpha_1, \alpha_2, \dots, \alpha_K$.

For completeness, we repeat the known result that if $L \sim \text{Poisson}(\lambda)$ and $x|L \sim \text{Multinomial}(L, D\theta)$ (which thus means that $L = \sum_m x_m$ with probability one), then x_1, x_2, \dots, x_M are mutually independent Poisson random variables with parameters $\lambda[D\theta]_1, \lambda[D\theta]_2, \dots, \lambda[D\theta]_M$. Indeed, we consider the following joint probability mass function where x and L are assumed to be non-negative integers:

$$\begin{aligned}
p(x, L|\theta, \lambda) &= p(L|\lambda)p(x|L, \theta) \\
&= \mathbb{1}_{\{L=\sum_m x_m\}} \frac{\exp(-\lambda)\lambda^L}{L!} \frac{L!}{\prod_m x_m!} \prod_m [D\theta]_m^{x_m} \\
&= \mathbb{1}_{\{L=\sum_m x_m\}} \exp(-\lambda \sum_m [D\theta]_m) \lambda^{\sum_m x_m} \prod_m \frac{[D\theta]_m^{x_m}}{x_m!} \\
&= \mathbb{1}_{\{L=\sum_m x_m\}} \prod_m \frac{\exp(-\lambda [D\theta]_m) (\lambda [D\theta]_m)^{x_m}}{x_m!} \\
&= \mathbb{1}_{\{L=\sum_m x_m\}} \prod_m \text{Poisson}(x_m; \lambda [D\theta]_m),
\end{aligned}$$

where in the third equation we used the fact that

$$\sum_m [D\theta]_m = \sum_{m,k} D_{mk} \theta_k = \sum_k \theta_k \sum_m D_{mk} = 1.$$

We thus have $p(x, L|\theta, \lambda) = p(L|x) \prod_m p(x_m|\lambda[D\theta]_m)$ where $p(L|x)$ is simply the deterministic distribution $\mathbb{1}_{\{L=\sum_m x_m\}}$ and $p(x_m|\lambda[D\theta]_m)$ for $m = 1, \dots, M$ are independent Poisson($\lambda[D\theta]_m$) distributions (and thus do not depend on L). Note that in the notation introduced in the paper, $D_{mk} = d_{km}$. Hence, by using the construction of the Dirichlet distribution from the normalization of independent gamma random variables, we can show that the LDA model with a gamma-Poisson prior over the length is equivalent to the following model (recall, that $c_0 = \sum_k c_k$):

$$\begin{aligned}
\lambda &\sim \text{Gamma}(c_0, b), \\
\theta &\sim \text{Dirichlet}(c), \\
x_m|\lambda, \theta &\sim \text{Poisson}([D(\lambda\theta)]_m).
\end{aligned} \tag{20}$$

More specifically, we complete the second part of the argument with the following properties. When $\alpha_1, \alpha_2, \dots, \alpha_K$ are mutually independent gamma random variables, each $\alpha_k \sim \text{Gamma}(c_k, b)$, their sum is also a gamma random variable $\sum_k \alpha_k \sim \text{Gamma}(\sum_k c_k, b)$. The former is equivalent to λ . It is known (e.g. [33]) that a Dirichlet random variable can be sampled by first sampling independent gamma random variables (α_k) and then dividing each of them by their sum (λ): $\theta_k = \alpha_k / \sum_{k'} \alpha_{k'}$, and, in other direction, the variables $\alpha_k = \lambda\theta_k$ are mutually independent, giving back the GP model (4).

A.2 The expectation and the variance of the document length for the GP model

From the derivations in Appendix A.1, it follows that the document length of the GP model (4) is a gamma-Poisson random variable, i.e. $L|\lambda \sim \text{Poisson}(\lambda)$ and $\lambda \sim \text{Gamma}(c_0, b)$. Therefore, the following follows from the law of total expectation and the law of total variance

$$\begin{aligned}
\mathbb{E}(L) &= \mathbb{E}[\mathbb{E}(L|\lambda)] = \mathbb{E}(\lambda) = c_0/b \\
\text{var}(L) &= \text{var}[\mathbb{E}(L|\lambda)] + \mathbb{E}[\text{var}(L|\lambda)] = \text{var}(\lambda) + \mathbb{E}(\lambda) = c_0/b + c_0/b^2
\end{aligned}$$

The first expression shows that the parameter b controls the expected document length $\mathbb{E}(L)$ for a given parameter c_0 : the smaller b , the larger $\mathbb{E}(L)$. On the other hand, if we allow c_0 to vary as well, only the ratio c_0/b is important for the document length. We can then interpret the role of c_0 as actually controlling the concentration of the distribution for the length L (through the variance). More specifically, we have that:

$$\frac{\text{var}(L)}{(\mathbb{E}(L))^2} = \frac{1}{\mathbb{E}(L)} + \frac{1}{c_0}. \quad (21)$$

For a fixed target document length $\mathbb{E}(L)$, we can increase the variance (and thus decrease the concentration) by using a smaller c_0 .

B Appendix. The cumulants of the GP and DICA models

B.1 Cumulants

For a random vector $x \in \mathbb{R}^M$, the first three cumulant tensors⁹ are

$$\begin{aligned} \text{cum}(x_m) &= \mathbb{E}(x_m), \\ \text{cum}(x_{m_1}, x_{m_2}) &= \mathbb{E}[(x_{m_1} - \mathbb{E}(x_{m_1}))(x_{m_2} - \mathbb{E}(x_{m_2}))] = \text{cov}(x_{m_1}, x_{m_2}), \\ \text{cum}(x_{m_1}, x_{m_2}, x_{m_3}) &= \mathbb{E}[(x_{m_1} - \mathbb{E}(x_{m_1}))(x_{m_2} - \mathbb{E}(x_{m_2}))(x_{m_3} - \mathbb{E}(x_{m_3}))]. \end{aligned}$$

Note that the 2nd and 3rd cumulants coincide with the 2nd and 3rd central moments (but not for higher orders). In the following, $\text{cum}(x, x, x) \in \mathbb{R}^{M \times M \times M}$ denotes the third order tensor with elements $\text{cum}(x_{m_1}, x_{m_2}, x_{m_3})$. Some of the properties of cumulants are listed below (see [19, chap. 5]). The most important property that motivate us to use cumulants in this paper (and the ICA literature) is the **independence** property, which says that the cumulant tensor for a random vector with independent components is diagonal (this property *does not* hold for the (non-central) moment tensors of any order, and neither for the central moments of order 4 or more).

- **Independence.** If the elements of $x \in \mathbb{R}^M$ are independent, then their cross-cumulants are zero as soon as two indices are different, i.e. $\text{cum}(x_{m_1}, x_{m_2}) = \delta(m_1, m_2)\mathbb{E}[(x_{m_1} - \mathbb{E}(x_{m_1}))^2]$ and $\text{cum}(x_{m_1}, x_{m_2}, x_{m_3}) = \delta(m_1, m_2, m_3)\mathbb{E}[(x_{m_1} - \mathbb{E}(x_{m_1}))^3]$, where δ is the Kronecker delta.
- **Multilinearity.** If two random vectors $y \in \mathbb{R}^M$ and $\alpha \in \mathbb{R}^K$ are linearly dependent, i.e. $y = D\alpha$ for some $D \in \mathbb{R}^{M \times K}$, then

$$\begin{aligned} \text{cum}(y_m) &= \sum_k \text{cum}(\alpha_k) D_{mk}, \\ \text{cum}(y_{m_1}, y_{m_2}) &= \sum_{k_1, k_2} \text{cum}(\alpha_{k_1}, \alpha_{k_2}) D_{m_1 k_1} D_{m_2 k_2}, \\ \text{cum}(y_{m_1}, y_{m_2}, y_{m_3}) &= \sum_{k_1, k_2, k_3} \text{cum}(\alpha_{k_1}, \alpha_{k_2}, \alpha_{k_3}) D_{m_1 k_1} D_{m_2 k_2} D_{m_3 k_3}, \end{aligned}$$

which can also be denoted¹⁰ by

$$\begin{aligned} \mathbb{E}(y) &= D\mathbb{E}(\alpha), \\ \text{cov}(y, y) &= D\text{cov}(\alpha, \alpha)D^\top, \\ \text{cum}(y, y, y) &= \text{cum}(\alpha, \alpha, \alpha)(D^\top, D^\top, D^\top). \end{aligned}$$

- **The law of total cumulance.** For two random vectors $x \in \mathbb{R}^M$ and $y \in \mathbb{R}^M$, it holds

$$\begin{aligned} \text{cum}(x_m) &= \mathbb{E}[\mathbb{E}(x_m|y)], \\ \text{cum}(x_{m_1}, x_{m_2}) &= \mathbb{E}[\text{cov}(x_{m_1}, x_{m_2}|y)] + \text{cov}[\mathbb{E}(x_{m_1}|y), \mathbb{E}(x_{m_2}|y)], \\ \text{cum}(x_{m_1}, x_{m_2}, x_{m_3}) &= \mathbb{E}[\text{cum}(x_{m_1}, x_{m_2}, x_{m_3}|y)] + \text{cum}[\mathbb{E}(x_{m_1}|y), \mathbb{E}(x_{m_2}|y), \mathbb{E}(x_{m_3}|y)] \\ &\quad + \text{cov}[\mathbb{E}(x_{m_1}|y), \text{cov}(x_{m_2}, x_{m_3}|y)] \\ &\quad + \text{cov}[\mathbb{E}(x_{m_2}|y), \text{cov}(x_{m_1}, x_{m_3}|y)] \\ &\quad + \text{cov}[\mathbb{E}(x_{m_3}|y), \text{cov}(x_{m_1}, x_{m_2}|y)]. \end{aligned}$$

⁹Strictly speaking, the (scalar) n -th cumulant κ_n of a random variable X is defined via the cumulant-generating function $g(t)$, which is the natural logarithm of the moment-generating function, i.e. $g(t) := \log \mathbb{E}[e^{tX}]$. The cumulant κ_n is then obtained from a power series expansion of the cumulant-generating function, that is $g(t) = \sum_{n=1}^{\infty} \kappa_n t^n / n!$ [Wikipedia].

¹⁰In [3], given a tensor $T \in \mathbb{R}^{K \times K \times K}$, $T(D^\top, D^\top, D^\top)$ is referred to as the multilinear map. In [34], the same entity is denoted by $T \times_1 D^\top \times_2 D^\top \times_3 D^\top$, where \times_n denotes the n -mode tensor-matrix product.

Note that the first expression is also well known as the law of total expectation or the tower property, while the second one is known as the law of total covariance.

B.2 The third cumulant of the GP/DICA models

In this section, by analogy with Section 3.1, we derive the third GP/DICA cumulant.

As the third cumulant of a Poisson random variable x_m with parameter y_m is $\mathbb{E}((x_m - \mathbb{E}(x_m))^3 | y_m) = y_m$, then by the independence property of cumulants from Section B.1, the cumulant of $x|y$ is diagonal:

$$\text{cum}(x_{m_1}, x_{m_2}, x_{m_3} | y) = \delta(m_1, m_2, m_3) y_{m_1}.$$

Substituting the cumulant of $x|y$ into the law of total cumulance, we obtain

$$\begin{aligned} \text{cum}(x_{m_1}, x_{m_2}, x_{m_3}) &= \mathbb{E}[\text{cum}(x_{m_1}, x_{m_2}, x_{m_3} | y)] \\ &\quad + \text{cum}[\mathbb{E}(x_{m_1} | y), \mathbb{E}(x_{m_2} | y), \mathbb{E}(x_{m_3} | y)] + \text{cov}[\mathbb{E}(x_{m_1} | y), \text{cov}(x_{m_2}, x_{m_3} | y)] \\ &\quad + \text{cov}[\mathbb{E}(x_{m_2} | y), \text{cov}(x_{m_1}, x_{m_3} | y)] + \text{cov}[\mathbb{E}(x_{m_3} | y), \text{cov}(x_{m_1}, x_{m_2} | y)] \\ &= \delta(m_1, m_2, m_3) \mathbb{E}(y_{m_1}) + \text{cum}(y_{m_1}, y_{m_2}, y_{m_3}) \\ &\quad + \delta(m_2, m_3) \text{cov}(y_{m_1}, y_{m_2}) + \delta(m_1, m_3) \text{cov}(y_{m_1}, y_{m_2}) + \delta(m_1, m_2) \text{cov}(y_{m_1}, y_{m_3}) \\ &= \delta(m_1, m_2, m_3) \mathbb{E}(x_{m_1}) + \text{cum}(y_{m_1}, y_{m_2}, y_{m_3}) \\ &\quad + \delta(m_2, m_3) \text{cov}(x_{m_1}, x_{m_2}) - \delta(m_1, m_2, m_3) \mathbb{E}(x_{m_1}) \\ &\quad + \delta(m_1, m_3) \text{cov}(x_{m_1}, x_{m_2}) - \delta(m_1, m_2, m_3) \mathbb{E}(x_{m_1}) \\ &\quad + \delta(m_1, m_2) \text{cov}(x_{m_1}, x_{m_3}) - \delta(m_1, m_2, m_3) \mathbb{E}(x_{m_1}) \\ &= \text{cum}(y_{m_1}, y_{m_2}, y_{m_3}) - 2\delta(m_1, m_2, m_3) \mathbb{E}(x_{m_1}) \\ &\quad + \delta(m_2, m_3) \text{cov}(x_{m_1}, x_{m_2}) + \delta(m_1, m_3) \text{cov}(x_{m_1}, x_{m_2}) + \delta(m_1, m_2) \text{cov}(x_{m_1}, x_{m_3}) \\ &= [\text{cum}(\alpha, \alpha, \alpha)(D^\top, D^\top, D^\top)]_{m_1 m_2 m_3} - 2\delta(m_1, m_2, m_3) \mathbb{E}(x_{m_1}) \\ &\quad + \delta(m_2, m_3) \text{cov}(x_{m_1}, x_{m_2}) + \delta(m_1, m_3) \text{cov}(x_{m_1}, x_{m_2}) + \delta(m_1, m_2) \text{cov}(x_{m_1}, x_{m_3}), \end{aligned} \quad (22)$$

where, in the third equality, we used the previous result from (10) that $\text{cov}(y, y) = \text{cov}(x, x) - \text{diag}(\mathbb{E}(x))$.

B.3 The diagonal structure of the GP/DICA cumulants

In this section, we provide detailed derivation of the diagonal structure (12) of the matrix S (11) and the diagonal structure (14) of the tensor T (13).

From the independence of $\alpha_1, \alpha_2, \dots, \alpha_K$ and by the independence property of cumulants from Section B.1, it follows that $\text{cov}(\alpha, \alpha)$ is a diagonal matrix and $\text{cum}(\alpha, \alpha, \alpha)$ is a diagonal tensor, i.e. $\text{cov}(\alpha_{k_1}, \alpha_{k_2}) = \delta(k_1, k_2) \text{cov}(\alpha_{k_1}, \alpha_{k_2})$ and $\text{cum}(\alpha_{k_1}, \alpha_{k_2}, \alpha_{k_3}) = \delta(k_1, k_2, k_3) \text{cum}(\alpha_{k_1}, \alpha_{k_1}, \alpha_{k_1})$. Therefore, the following holds

$$\begin{aligned} \text{cov}(y_{m_1}, y_{m_2}) &= \sum_k \text{cov}(\alpha_k, \alpha_k) D_{m_1 k} D_{m_2 k}, \\ \text{cum}(y_{m_1}, y_{m_2}, y_{m_3}) &= \sum_k \text{cum}(\alpha_k, \alpha_k, \alpha_k) D_{m_1 k} D_{m_2 k} D_{m_3 k}, \end{aligned}$$

which we can rewrite in a matrix/tensor form as

$$\begin{aligned} \text{cov}(y, y) &= \sum_k \text{cov}(\alpha_k, \alpha_k) d_k d_k^\top, \\ \text{cum}(y, y, y) &= \sum_k \text{cum}(\alpha_k, \alpha_k, \alpha_k) d_k \otimes d_k \otimes d_k. \end{aligned}$$

Moving $\text{cov}(y, y) / \text{cum}(y, y, y)$ in the expression for $\text{cov}(x, x)$ (10) / $\text{cum}(x, x, x)$ (22) on one side of equality and all other terms on the other side, we define matrix $S \in \mathbb{R}^{M \times M}$ / tensor $T \in \mathbb{R}^{M \times M \times M}$

as follows

$$S := \text{cov}(x, x) - \text{diag}(\mathbb{E}(x)), \quad (23)$$

$$\begin{aligned} T_{m_1 m_2 m_3} &:= \text{cum}(x_{m_1}, x_{m_2}, x_{m_3}) + 2\delta(m_1, m_2, m_3)\mathbb{E}(x_{m_1}) \\ &\quad - \delta(m_2, m_3)\text{cov}(x_{m_1}, x_{m_2}) \\ &\quad - \delta(m_1, m_3)\text{cov}(x_{m_1}, x_{m_2}) \\ &\quad - \delta(m_1, m_2)\text{cov}(x_{m_1}, x_{m_3}). \end{aligned} \quad (24)$$

By construction, $S = \text{cov}(y, y)$ and $T = \text{cum}(y, y, y)$ and, therefore, it holds that

$$S = \sum_k \text{cov}(\alpha_k, \alpha_k) d_k d_k^\top, \quad (25)$$

$$T = \sum_k \text{cum}(\alpha_k, \alpha_k, \alpha_k) d_k \otimes d_k \otimes d_k. \quad (26)$$

This means that both the matrix S and the tensor T are sums of rank-1 matrices and tensors, respectively¹¹. This structure of the matrix S and the tensor T is the basis for the algorithms considered in this paper.

B.4 Unbiased finite sample estimators for the GP/DICA cumulants

Given a sample $\{x_1, x_2, \dots, x_N\}$, we obtain a finite sample estimate \hat{S} of S (11) / \hat{T} of T (13) for the GP/DICA cumulants:

$$\hat{S} := \widehat{\text{cov}}(x, x) - \text{diag}(\hat{\mathbb{E}}(x)), \quad (27)$$

$$\begin{aligned} \hat{T}_{m_1 m_2 m_3} &:= \widehat{\text{cum}}(x_{m_1}, x_{m_2}, x_{m_3}) + 2\delta(m_1, m_2, m_3)\hat{\mathbb{E}}(x_{m_1}) \\ &\quad - \delta(m_2, m_3)\widehat{\text{cov}}(x_{m_1}, x_{m_2}) \\ &\quad - \delta(m_1, m_3)\widehat{\text{cov}}(x_{m_1}, x_{m_2}) \\ &\quad - \delta(m_1, m_2)\widehat{\text{cov}}(x_{m_1}, x_{m_3}), \end{aligned} \quad (28)$$

where unbiased estimators of the first three cumulants are

$$\begin{aligned} \hat{\mathbb{E}}(x_{m_1}) &= \frac{1}{N} \sum_n x_{nm_1}, \\ \widehat{\text{cov}}(x_{m_1}, x_{m_2}) &= \frac{1}{N-1} \sum_n z_{nm_1} z_{nm_2}, \\ \widehat{\text{cum}}(x_{m_1}, x_{m_2}, x_{m_3}) &= \frac{N}{(N-1)(N-2)} \sum_n z_{nm_1} z_{nm_2} z_{nm_3}, \end{aligned} \quad (29)$$

where the word vocabulary indexes are $m_1, m_2, m_3 = 1, 2, \dots, M$ and the centered documents $z_{nm} := x_{nm} - \hat{\mathbb{E}}(x_m)$. (The latter is introduced only for compact representation of (29) and is different from z in the LDA model.)

C Appendix. The sketch of the proof for Proposition 3.1

C.1 Expected squared error for the sample expectation

The sample expectation is $\hat{\mathbb{E}}(x) = \frac{1}{N} \sum_n x_n$ is an unbiased estimator of the expectation and:

$$\begin{aligned} \mathbb{E} \left(\|\hat{\mathbb{E}}(x) - \mathbb{E}(x)\|_2^2 \right) &= \sum_m \mathbb{E} \left[\left(\hat{\mathbb{E}}(x_m) - \mathbb{E}(x_m) \right)^2 \right] \\ &= \frac{1}{N^2} \sum_m \left[\mathbb{E} \left(\sum_n (x_{nm} - \mathbb{E}(x_m))^2 \right) + \mathbb{E} \left(\sum_n \sum_{n' \neq n} (x_{nm} - \mathbb{E}(x_m)) (x_{n'm} - \mathbb{E}(x_m)) \right) \right] \\ &= \frac{1}{N} \sum_m \mathbb{E} \left[(x_m - \mathbb{E}(x_m))^2 \right] = \frac{1}{N} \sum_m \text{var}(x_m). \end{aligned}$$

¹¹For tensors, such decomposition is also known under the names CANDECOMP/PARAFAC or, simply, the CP decomposition (see, e.g., [34]).

Further, by the law of total variance:

$$\begin{aligned}\mathbb{E}\left(\|\widehat{\mathbb{E}}(x) - \mathbb{E}(x)\|_2^2\right) &= \frac{1}{N} \sum_m [\mathbb{E}(\text{var}(x_m|y)) + \text{var}(\mathbb{E}(x_m|y))] = \frac{1}{N} \sum_m [\mathbb{E}(y_m) + \text{var}(y_m)] \\ &= \frac{1}{N} \left[\sum_k \mathbb{E}(\alpha_k) + \sum_k \langle d_k, d_k \rangle \text{var}(\alpha_k) \right],\end{aligned}$$

using the fact that $\sum_m D_{mk} = 1$ for any k .

C.2 Expected squared error for the sample covariance

The following finite sample estimator of the covariance $\text{cov}(x, x) = \mathbb{E}(xx^\top) - \mathbb{E}(x)\mathbb{E}(x)^\top$

$$\begin{aligned}\widehat{\text{cov}}(x, x) &= \frac{1}{N-1} \sum_n x_n x_n^\top - \widehat{\mathbb{E}}(x) \widehat{\mathbb{E}}(x)^\top = \frac{1}{N-1} \sum_n \left(x_n x_n^\top - \frac{1}{N^2} \sum_{n'} \sum_{n''} x_{n'} x_{n''}^\top \right) \\ &= \frac{1}{N} \sum_n \left(x_n x_n^\top - \frac{1}{N-1} x_n \sum_{n' \neq n} x_{n'}^\top \right)\end{aligned}\tag{30}$$

is unbiased, i.e. $\mathbb{E}(\widehat{\text{cov}}(x, x)) = \text{cov}(x, x)$. Its squared error is

$$\mathbb{E}(\|\widehat{\text{cov}}(x, x) - \text{cov}(x, x)\|_F^2) = \sum_{m, m'} \mathbb{E} \left[(\widehat{\text{cov}}(x_m, x_{m'}) - \mathbb{E}[\widehat{\text{cov}}(x_m, x_{m'})])^2 \right].$$

The m, m' -th element of the sum above is equal to

$$\begin{aligned}& \frac{1}{N^2} \sum_{n, n'} \text{cov} \left(x_{nm} x_{nm'} - \frac{1}{N-1} x_{nm} \sum_{n'' \neq n} x_{n''m'}, \quad x_{n'm} x_{n'm'} - \frac{1}{N-1} x_{n'm} \sum_{n'' \neq n'} x_{n''m'} \right) \\ &= \frac{1}{N^2} \sum_{n, n'} \text{cov}(x_{nm} x_{nm'}, x_{n'm} x_{n'm'}) - \frac{2}{N^2(N-1)} \sum_{n, n'} \text{cov} \left(x_{nm} \sum_{n'' \neq n} x_{n''m'}, x_{n'm} x_{n'm'} \right) \\ &+ \frac{1}{N^2(N-1)^2} \sum_{n, n'} \text{cov} \left(x_{nm} \sum_{n'' \neq n} x_{n''m'}, x_{n'm} \sum_{n'' \neq n'} x_{n''m'} \right) \\ &= \frac{1}{N^2} \sum_n \text{cov}(x_{nm} x_{nm'}, x_{nm} x_{nm'}) \\ &- \frac{2}{N^2(N-1)} \left[\sum_n \sum_{n'' \neq n} \text{cov}(x_{nm} x_{n''m'}, x_{nm} x_{nm'}) + \sum_n \sum_{n' \neq n} \text{cov}(x_{nm} x_{n'm'}, x_{n'm} x_{n'm'}) \right] \\ &+ \frac{1}{N^2(N-1)^2} \left[\sum_n \sum_{n'' \neq n} \sum_{n''' \neq n} \text{cov}(x_{nm} x_{n''m'}, x_{nm} x_{n'''m'}) + \sum_{n'} \sum_{n \neq n'} \sum_{n'' \neq n'} \text{cov}(x_{nm} x_{n''m'}, x_{n'm} x_{nm'}) \right] \\ &+ \frac{1}{N^2(N-1)^2} \left[\sum_{n'} \sum_{n \neq n'} \sum_{n'' \neq n'} \text{cov}(x_{nm} x_{n'm'}, x_{n'm} x_{n''m'}) + \sum_{n'} \sum_{n \neq n'} \sum_{n'' \neq n'} \text{cov}(x_{nm} x_{n''m'}, x_{n'm} x_{n''m'}) \right],\end{aligned}$$

where we used mutual independence of the observations x_n in a sample $\{x_n\}_{n=1}^N$ to conclude that the covariance between two expressions involving only independent variables is zero. Further:

$$\begin{aligned}\mathbb{E}(\|\widehat{\text{cov}}(x, x) - \text{cov}(x, x)\|_F^2) &= \frac{1}{N^2} \sum_{m, m'} N \left(\mathbb{E}(x_m^2 x_{m'}^2) - [\mathbb{E}(x_m x_{m'})]^2 \right) \\ &- \frac{4}{N^2(N-1)} \sum_{m, m'} N(N-1) \left(\mathbb{E}(x_m^2 x_{m'}) \mathbb{E}(x_{m'}) - \mathbb{E}(x_m x_{m'}) \mathbb{E}(x_m) \mathbb{E}(x_{m'}) \right) \\ &+ \frac{2}{N^2(N-1)^2} \sum_{m, m'} N(N-1)(N-2) \left(\mathbb{E}(x_m^2) [\mathbb{E}(x_{m'})]^2 - [\mathbb{E}(x_m)]^2 [\mathbb{E}(x_{m'})]^2 \right) \\ &+ \frac{2}{N^2(N-1)^2} \sum_{m, m'} N(N-1)(N-2) \left(\mathbb{E}(x_m x_{m'}) \mathbb{E}(x_m) \mathbb{E}(x_{m'}) - [\mathbb{E}(x_m)]^2 [\mathbb{E}(x_{m'})]^2 \right) + O\left(\frac{1}{N^2}\right),\end{aligned}$$

which after simplification gives

$$\begin{aligned}\mathbb{E}(\|\widehat{\text{cov}}(x, x) - \text{cov}(x, x)\|_F^2) &= \frac{1}{N} \sum_{m, m'} \left[\text{var}(x_m x_{m'}) + 2 [\mathbb{E}(x_m)]^2 \text{var}(x_{m'}) \right] \\ &+ \frac{1}{N} \sum_{m, m'} [2\mathbb{E}(x_m)\mathbb{E}(x_{m'})\text{cov}(x_m, x_{m'}) - 4\mathbb{E}(x_m)\text{cov}(x_m x_{m'}, x_{m'})] + O\left(\frac{1}{N^2}\right),\end{aligned}$$

where in the last equality, by symmetry, the summation indexes m and m' can be exchanged. As $x_m \sim \text{Poisson}(y_m)$, by the law of total expectation and law of total covariance, it follows, for $m \neq m'$ (and using the auxiliary expressions from Section C.4):

$$\begin{aligned}\text{var}(x_m x_{m'}) &= \mathbb{E}(x_m^2 x_{m'}^2) - [\mathbb{E}(x_m x_{m'})]^2 \\ &= \mathbb{E}[y_m^2 y_{m'}^2 + y_m^2 y_{m'} + y_m y_{m'}^2 + y_m y_{m'}] - [\mathbb{E}(y_m y_{m'})]^2 \\ [\mathbb{E}(x_m)]^2 \text{var}(x_{m'}) &= [\mathbb{E}(y_m)]^2 [\mathbb{E}[\text{var}(x'_{m'}|y)] + \text{var}[\mathbb{E}(x'_{m'}|y)]] \\ &= [\mathbb{E}(y_m)]^2 \mathbb{E}(y_{m'}) + [\mathbb{E}(y_m)]^2 \mathbb{E}(y_{m'}^2) - [\mathbb{E}(y_m)]^2 [\mathbb{E}(y_{m'})]^2, \\ \mathbb{E}(x_m)\mathbb{E}(x_{m'})\text{cov}(x_m, x_{m'}) &= \mathbb{E}(y_m y_{m'})\mathbb{E}(y_m)\mathbb{E}(y_{m'}) - [\mathbb{E}(y_m)]^2 [\mathbb{E}(y_{m'})]^2, \\ \mathbb{E}(x_m)\text{cov}(x_m x_{m'}, x_{m'}) &= \mathbb{E}(y_m) [\mathbb{E}(y_m y_{m'}) + \mathbb{E}(y_m y_{m'}^2) - \mathbb{E}(y_m y_{m'})\mathbb{E}(y_{m'})].\end{aligned}$$

Now, considering the $m = m'$ case, we have:

$$\begin{aligned}\text{var}(x_m^2) &= \mathbb{E}[\mathbb{E}(x_m^4|y)] - (\mathbb{E}[\mathbb{E}(x_m^2|y)])^2 \\ &= \mathbb{E}[y_m^4 + 6y_m^3 + 7y_m^2 + y_m] - (\mathbb{E}[y_m^2 + y_m])^2 \\ \mathbb{E}(x_m)\mathbb{E}(x_m)\text{cov}(x_m, x_m) &= \mathbb{E}(y_m^2) [\mathbb{E}(y_m^2) + \mathbb{E}(y_m) - [\mathbb{E}(y_m)]^2], \\ \mathbb{E}(x_m)\text{cov}(x_m^2, x_m) &= \mathbb{E}(y_m) [\mathbb{E}(y_m^3) + 3\mathbb{E}(y_m^2) + \mathbb{E}(y_m) - \mathbb{E}(y_m) [\mathbb{E}(y_m^2) + \mathbb{E}(y_m)]].\end{aligned}$$

Substitution of $y_m = \sum_k D_{mk} \alpha_k$ gives the following

$$\begin{aligned}\mathbb{E}(\|\widehat{\text{cov}}(x, x) - \text{cov}(x, x)\|_F^2) &= \frac{1}{N} \sum_{k, k', k'', k'''} \langle d_k, d_{k'} \rangle \langle d_{k''}, d_{k'''} \rangle \mathcal{A}_{kk'k''k'''} \\ &+ \frac{1}{N} \sum_{k, k', k''} \langle d_k, d_{k'} \rangle \langle d_{k''}, \vec{1} \rangle \mathcal{B}_{kk'k''} + \frac{1}{N} \sum_{k, k'} \langle d_k, \vec{1} \rangle \langle d_{k'}, \vec{1} \rangle \mathbb{E}(\alpha_k \alpha_{k'}) + O\left(\frac{1}{N^2}\right),\end{aligned}$$

where $\vec{1}$ is the vector with all the elements equal to 1 and

$$\begin{aligned}\mathcal{A}_{kk'k''k'''} &= \mathbb{E}(\alpha_k \alpha_{k'} \alpha_{k''} \alpha_{k'''}) - \mathbb{E}(\alpha_k \alpha_{k'}) \mathbb{E}(\alpha_{k''} \alpha_{k'''}) + 2\mathbb{E}(\alpha_k) \mathbb{E}(\alpha_{k'}) \mathbb{E}(\alpha_{k''} \alpha_{k'''}) \\ &- 2\mathbb{E}(\alpha_k) \mathbb{E}(\alpha_{k'}) \mathbb{E}(\alpha_{k''}) \mathbb{E}(\alpha_{k'''}) + 2\mathbb{E}(\alpha_k \alpha_{k''}) \mathbb{E}(\alpha_{k'}) \mathbb{E}(\alpha_{k'''}) - 2\mathbb{E}(\alpha_k) \mathbb{E}(\alpha_{k'}) \mathbb{E}(\alpha_{k''}) \mathbb{E}(\alpha_{k'''}) \\ &- 4\mathbb{E}(\alpha_k) \mathbb{E}(\alpha_{k'} \alpha_{k''} \alpha_{k'''}) + 4\mathbb{E}(\alpha_k) \mathbb{E}(\alpha_{k'} \alpha_{k''}) \mathbb{E}(\alpha_{k'''}), \\ \mathcal{B}_{kk'k''} &= 2\mathbb{E}(\alpha_k \alpha_{k'} \alpha_{k''}) + 4\mathbb{E}(\alpha_k) \mathbb{E}(\alpha_{k'}) \mathbb{E}(\alpha_{k''}) - 4\mathbb{E}(\alpha_k) \mathbb{E}(\alpha_{k'} \alpha_{k''}),\end{aligned}$$

where we used the expressions from Section C.4.

C.3 Expected squared error of the estimator \widehat{S} for the GP/DICA cumulants

As the estimator \widehat{S} (27) of S (11) is unbiased, its expected squared error is

$$\begin{aligned}\mathbb{E}[\|\widehat{S} - S\|_F^2] &= \mathbb{E}\left[\left\|\widehat{\text{cov}}(x, x) - \text{cov}(x, x) + \left(\text{diag}[\widehat{\mathbb{E}}(x)] - \text{diag}[\mathbb{E}(x)]\right)\right\|_F^2\right] \\ &= \mathbb{E}\left[\|\widehat{\mathbb{E}}(x) - \mathbb{E}(x)\|_F^2\right] + \mathbb{E}\left[\|\widehat{\text{cov}}(x, x) - \text{cov}(x, x)\|_F^2\right] \\ &+ 2 \sum_m \mathbb{E}\left[\left(\widehat{\mathbb{E}}(x_m) - \mathbb{E}(x_m)\right) (\widehat{\text{cov}}(x_m, x_m) - \text{cov}(x_m, x_m))\right].\end{aligned}\tag{31}$$

As $\widehat{\mathbb{E}}(x_m)$ and $\widehat{\text{cov}}(x_m, x_m)$ are unbiased, the m -th element of the last sum is equal to

$$\begin{aligned}
& \text{cov} \left[\widehat{\mathbb{E}}(x_m), \widehat{\text{cov}}(x_m, x_m) \right] \\
&= \frac{1}{N^2} \sum_{n, n'} \text{cov} [x_{nm}, x_{n'm}^2] - \frac{1}{N^2(N-1)} \sum_{n, n', n'' \neq n'} \text{cov} [x_{nm}, x_{n'm} x_{n''m}] \\
&= \frac{1}{N^2} \sum_n \text{cov} [x_{nm}, x_{nm}^2] - \frac{2}{N^2(N-1)} \sum_{n, n' \neq n} \text{cov} [x_{nm}, x_{n'm} x_{nm}] + O\left(\frac{1}{N^2}\right) \\
&= \frac{1}{N} \mathbb{E}(x_m^3) - \frac{2}{N} \left(\mathbb{E}(x_m^2) \mathbb{E}(x_m) - [\mathbb{E}(x_m)]^3 \right) + O\left(\frac{1}{N^2}\right) \\
&\leq \frac{1}{N} \mathbb{E}(x_m^3) + \frac{2}{N} [\mathbb{E}(x_m)]^3 + O\left(\frac{1}{N^2}\right) = \frac{1}{N} \left[\mathbb{E}(y_m^3) + 3\mathbb{E}(y_m^2) + \mathbb{E}(y_m) + 2[\mathbb{E}(y_m)]^3 \right] + O\left(\frac{1}{N^2}\right),
\end{aligned}$$

where we neglected the negative term $-\mathbb{E}(x_m^2)\mathbb{E}(x_m)$ for the inequality, and the last equality follows from the expressions in Section C.4. Further, the fact that $y_m = \sum_k D_{mk} \alpha_k$ gives

$$\begin{aligned}
\sum_m \text{cov} \left[\widehat{\mathbb{E}}(x_m), \widehat{\text{cov}}(x_m, x_m) \right] &= \frac{1}{N} \sum_{k, k', k''} \langle d_k \circ d_{k'}, d_{k''} \rangle \mathcal{C}_{kk'k''} \\
&\quad + \frac{3}{N} \sum_{k, k'} \langle d_k, d_{k'} \rangle \mathbb{E}(\alpha_k \alpha_{k'}) + \frac{1}{N} \sum_k \langle d_k, \vec{1} \rangle \mathbb{E}(\alpha_k) + O\left(\frac{1}{N^2}\right),
\end{aligned}$$

where \circ denotes the elementwise Hadamard product and

$$\mathcal{C}_{kk'k''} = \mathbb{E}(\alpha_k \alpha_{k'} \alpha_{k''}) + 2\mathbb{E}(\alpha_k) \mathbb{E}(\alpha_{k'}) \mathbb{E}(\alpha_{k''}).$$

Plugging this and the expressions for $\mathbb{E}(\|\widehat{\mathbb{E}}(x) - \mathbb{E}(x)\|_F^2)$ and $\mathbb{E}(\|\widehat{\text{cov}}(x, x) - \text{cov}(x, x)\|_F^2)$ from Sections C.1 and C.2, respectively, into (31) gives

$$\begin{aligned}
\mathbb{E} \left[\|\widehat{S} - S\|_F^2 \right] &= \frac{1}{N} \left[\sum_k \langle d_k, d_k \rangle \text{var}(\alpha_k) + \sum_k \mathbb{E}(\alpha_k) + \sum_{k, k', k'', k'''} \langle d_k, d_{k'} \rangle \langle d_{k''}, d_{k'''} \rangle \mathcal{A}_{kk'k''k'''} \right] + O\left(\frac{1}{N^2}\right) \\
&\quad + \frac{1}{N} \left[\sum_{k, k', k''} [\langle d_k, d_{k'} \rangle \mathcal{B}_{kk'k''} + 2\langle d_k \circ d_{k'}, d_{k''} \rangle \mathcal{C}_{kk'k''}] + \sum_{k, k'} (1 + 6\langle d_k, d_{k'} \rangle) \mathbb{E}(\alpha_k \alpha_{k'}) + 2 \sum_k \mathbb{E}(\alpha_k) \right],
\end{aligned}$$

where we used that, by the simplex constraint on the topics, $\langle d_k, \vec{1} \rangle = 1$ for all k . To analyze this expression in more details, let us now consider the GP model, i.e. $\alpha_k \sim \text{Gamma}(c_k, b)$:

$$\begin{aligned}
\sum_{k, k', k'', k'''} \mathcal{A}_{kk'k''k'''} &\leq \frac{30c_0^4 + 23c_0^3 + 14c_0^2 + 8c_0}{b^4}, \quad \text{and} \quad \sum_{k, k', k''} \mathcal{B}_{kk'k''} \leq \frac{14c_0^3 + 12c_0^2 + 4c_0}{b^3}, \\
\sum_{k, k', k''} \mathcal{C}_{kk'k''} &\leq \frac{7c_0^3 + 6c_0^2 + 2c_0}{b^3}, \quad \text{and} \quad \sum_{k, k'} \mathbb{E}(\alpha_k \alpha_{k'}) \leq \frac{2c_0^2 + c_0}{b^2},
\end{aligned}$$

where we used the expressions from Section C.4, which gives

$$\begin{aligned}
\mathbb{E} \left[\|\widehat{S} - S\|_F^2 \right] &\leq \frac{\nu}{N} \left[\max_k \|d_k\|_2^2 \frac{c_0}{b^2} + \frac{c_0}{b} + \left(\max_{k, k'} \langle d_k, d_{k'} \rangle \right)^2 \max \left[\frac{c_0^4}{b^4}, \frac{c_0}{b^4} \right] + \max_{k, k'} \langle d_k, d_{k'} \rangle \max \left[\frac{c_0^3}{b^3}, \frac{c_0}{b^3} \right] \right] \\
&\quad + \frac{\nu}{N} \left[\left(\max_{k, k', k''} \langle d_k \circ d_{k'}, d_{k''} \rangle \right) \max \left[\frac{c_0^3}{b^3}, \frac{c_0}{b^3} \right] + \left(1 + \max_{k, k'} \langle d_k, d_{k'} \rangle \right) \max \left[\frac{c_0^2}{b^2}, \frac{c_0}{b^2} \right] \right] + O\left(\frac{1}{N^2}\right),
\end{aligned}$$

where $\nu \leq 30$ is a universal constant. As, by the Cauchy-Schwarz inequality, $\max_{k, k'} \langle d_k, d_{k'} \rangle \leq \max_k \|d_k\|_2^2 =: \Delta_1$ and $\max_{k, k', k''} \langle d_k \circ d_{k'}, d_{k''} \rangle \leq \max_k \|d_k\|_\infty \|d_k\|_2^2 \leq \max_k \|d_k\|_2^3 =: \Delta_2$ (note that for the topics in the simplex, $\Delta_2 \leq \Delta_1$ as well as $\Delta_1^2 \leq \Delta_2$), it follows that

$$\begin{aligned}
\mathbb{E} \left[\|\widehat{S} - S\|_F^2 \right] &\leq \frac{\nu}{N} \left[\Delta_1 \left(\frac{L^2}{\bar{c}_0} + \frac{L^3}{\bar{c}_0^2} \right) + L + \Delta_1^2 \frac{L^4}{\bar{c}_0^3} + \frac{L^2}{\bar{c}_0^2} + \Delta_2 \frac{L^3}{\bar{c}_0^2} \right] + O\left(\frac{1}{N^2}\right) \\
&\leq \frac{2\nu}{N} \frac{1}{\bar{c}_0^3} [\Delta_1^2 L^4 + \bar{c}_0 \Delta_1 L^3 + \bar{c}_0^2 L^2 + \bar{c}_0^3 L] + O\left(\frac{1}{N^2}\right),
\end{aligned}$$

where $\bar{c}_0 = \min(1, c_0) \leq 1$ and, from Section A.2, $c_0 = bL$ where L is the expected document length. The second term $\bar{c}_0 \Delta_1 L^3$ cannot be dominant as the system $\bar{c}_0 \Delta_1 L^3 > \bar{c}_0^2 L^2$ and $\bar{c}_0 \Delta_1 L^3 > \Delta_1^2 L^4$ is infeasible. Also, with the reasonable assumption that $L \geq 1$, we also have that the 4th term $\bar{c}_0^3 L \leq \bar{c}_0^2 L^2$. Therefore,

$$\mathbb{E} \left[\|\hat{S} - S\|_F^2 \right] \leq \frac{3\nu}{N} \max [\Delta_1^2 L^4, \bar{c}_0^2 L^2] + O \left(\frac{1}{N^2} \right).$$

C.4 Auxiliary expressions

As $\{x_m\}_{m=1}^M$ are conditionally independent given y in the DICA model (4), we have the following expressions by using the law of total expectation for $m \neq m'$ and using the moments of the Poisson distribution with parameter y_m :

$$\begin{aligned} \mathbb{E}(x_m) &= \mathbb{E}[\mathbb{E}(x_m|y_m)] = \mathbb{E}(y_m), \\ \mathbb{E}(x_m^2) &= \mathbb{E}[\mathbb{E}(x_m^2|y_m)] = \mathbb{E}(y_m^2) + \mathbb{E}(y_m), \\ \mathbb{E}(x_m^3) &= \mathbb{E}[\mathbb{E}(x_m^3|y_m)] = \mathbb{E}(y_m^3) + 3\mathbb{E}(y_m^2) + \mathbb{E}(y_m), \\ \mathbb{E}(x_m^4) &= \mathbb{E}[\mathbb{E}(x_m^4|y_m)] = \mathbb{E}(y_m^4) + 6\mathbb{E}(y_m^3) + 7\mathbb{E}(y_m^2) + \mathbb{E}(y_m), \\ \mathbb{E}(x_m x_{m'}) &= \mathbb{E}[\mathbb{E}(x_m x_{m'}|y)] = \mathbb{E}[\mathbb{E}(x_m|y_m)\mathbb{E}(x_{m'}|y_{m'})] = \mathbb{E}(y_m y_{m'}), \\ \mathbb{E}(x_m x_{m'}^2) &= \mathbb{E}[\mathbb{E}(x_m x_{m'}^2|y)] = \mathbb{E}[\mathbb{E}(x_m|y_m)\mathbb{E}(x_{m'}^2|y_{m'})] = \mathbb{E}(y_m y_{m'}^2) + \mathbb{E}(y_m y_{m'}), \\ \mathbb{E}(x_m^2 x_{m'}^2) &= \mathbb{E}[\mathbb{E}(x_m^2|y_m)\mathbb{E}(x_{m'}^2|y_{m'})] = \mathbb{E}(y_m^2 y_{m'}^2) + \mathbb{E}(y_m^2 y_{m'}) + \mathbb{E}(y_m y_{m'}^2) + \mathbb{E}(y_m y_{m'}). \end{aligned}$$

Moreover, the moments of $\alpha_k \sim \text{Gamma}(c_k, b)$ are

$$\mathbb{E}(\alpha_k) = \frac{c_k}{b}, \quad \mathbb{E}(\alpha_k^2) = \frac{c_k^2 + c_k}{b^2}, \quad \mathbb{E}(\alpha_k^3) = \frac{c_k^3 + 3c_k^2 + 2c_k}{b^3}, \quad \mathbb{E}(\alpha_k^4) = \frac{c_k^4 + 6c_k^3 + 11c_k^2 + 6c_k}{b^4}, \quad \text{etc.}$$

C.5 Analysis of whitening and recovery error

We can follow a similar analysis as in Appendix C of [2] to derive the topic recovery error given the sample estimate error. In particular, if we define the following sampling errors E_S and E_T :

$$\begin{aligned} \|\hat{S} - S\| &\leq E_S, \\ \|\hat{T}(u) - T(u)\| &\leq \|u\|_2 E_T, \end{aligned}$$

then the following form of their Lemma C.2 holds for both the LDA moments and the GP/DICA cumulants:

$$\|\widehat{W}\widehat{T}(\widehat{W}^\top u)\widehat{W}^\top - WT(W^\top u)W^\top\| \leq \nu \left[\frac{(\max_k \gamma_k)E_S}{\sigma_K(\tilde{D})^2} + \frac{E_T}{\sigma_K(\tilde{D})^3} \right], \quad (32)$$

where $\sigma_k(\cdot)$ denotes the k -th singular value of a matrix, ν is some universal constant, and in both cases \tilde{D} was defined such that $S = \tilde{D}\tilde{D}^\top$. For the LDA moments, $\gamma_k = 2\sqrt{\frac{c_0(c_0+1)}{c_k(c_0+2)^2}}$, whereas for the GP/DICA cumulants, γ_k takes the simpler form $\gamma_k := \text{cum}(\alpha_k)/[\text{var}(\alpha_k)]^{3/2} = 2/\sqrt{c_k}$.

We note that the scaling for S is $O(L^2)$ for the GP/DICA cumulants, in contrast to $O(1)$ for the LDA moments. Thus, to compare the upper bound (32) for the two types of moments, we need to put it in quantities which are common. In the first section of the Appendix C of [2], it was mentioned that $\sigma_K(\tilde{D}) \geq \sqrt{\frac{c_{\min}}{c_0(c_0+1)}}\sigma_K(D)$ for the LDA moments, where $c_{\min} := \min_k c_k$. In contrast, for the GP/DICA cumulants, we can show that $\sigma_K(\tilde{D}) \geq L\sqrt{\frac{c_{\min}}{c_0}}\sigma_K(D)$, where $L := c_0/b$ is the average length of a document in the GP model. Using this lower bound for the singular vector, we thus get the following bound in the case of the GP cumulant:

$$\|\widehat{W}\widehat{T}(\widehat{W}^\top u)\widehat{W}^\top - WT(W^\top u)W^\top\| \leq \frac{\nu}{c_{\min}^{3/2}} \left[\frac{E_S}{L^2} \frac{2c_0^2}{[\sigma_K(D)]^2} + \frac{E_T}{L^3} \frac{c_0^3}{[\sigma_K(D)]^3} \right]. \quad (33)$$

The $c_{\min}^{3/2}$ factor is common for both the LDA moment and GP cumulant, but as we mentioned after Proposition 3.1, the sample error E_S term gets divided by L^2 for the GP cumulant, as expected.

The recovery error bound in [2] is based on the bound (33), and thus by showing that the error E_S/L^2 for the GP cumulant is lower than the E_S term for the LDA moment, we expect to also gain a similar gain for the recovery error, as the rest of the argument is the same for both types of moments (see Appendix C.2, C.3 and C.4 in [2] for the completion).

D Appendix. The LDA moments

D.1 Our notation

The LDA moments were derived in [1]. Note that the full version of the paper with proofs appeared in [2] and a later version of the paper also appeared in [32]. In this section, we recall the form of the LDA moments using our notation. This section does not contain any novel results and is included for the reader's convenience. We also refer to this section when deriving the practical expressions for computation of the sample estimates of the LDA moments in Appendix E.3.

For deriving the LDA moments, a document is assumed to be composed of at least three tokens: $L \geq 3$. As the LDA generative model (1) is only defined *conditional* on the length L , this is not too problematic. But given that we present models in this paper which also model L , we mention for clarity that we can suppose that all expectations and probabilities defined below are implicitly conditioning on $L \geq 3$.¹² The theoretical LDA moments are derived only using the first three words w_1, w_2 and w_3 of a document. But note that since the words w_ℓ 's are conditionally i.i.d. given θ (for $1 \leq \ell \leq L$), we have $M_3 := \mathbb{E}(w_1 \otimes w_2 \otimes w_3) = \mathbb{E}(w_{\ell_1} \otimes w_{\ell_2} \otimes w_{\ell_3})$ for any three distinct tokens ℓ_1, ℓ_2 and ℓ_3 . The tensor M_3 is thus symmetric, and could have been defined using any distinct ℓ_1, ℓ_2 and ℓ_3 that are less than L . To highlight this arbitrary choice and to make the links with the U-statistics estimator presented later, we thus use generic distinct ℓ_1, ℓ_2 and ℓ_3 in the definition of the LDA moments below, instead of $\ell_1 = 1, \ell_2 = 2$ and $\ell_3 = 1$ as in [1].

Using this notation, then by the law of total expectation and the properties of the Dirichlet distribution, the non-central moments¹³ of the LDA model (1) take the form [1]:

$$M_1 = \mathbb{E}(w_{\ell_1}) = D \frac{c}{c_0}, \quad (34)$$

$$M_2 = \mathbb{E}(w_{\ell_1} w_{\ell_2}^\top) = \frac{c_0}{c_0 + 1} M_1 M_1^\top + \frac{1}{c_0(c_0 + 1)} D \text{diag}(c) D^\top, \quad (35)$$

$$\begin{aligned} M_3 &= \mathbb{E}(w_{\ell_1} \otimes w_{\ell_2} \otimes w_{\ell_3}) \\ &= \frac{c_0}{c_0 + 2} [\mathbb{E}(w_{\ell_1} \otimes w_{\ell_2} \otimes M_1) + \mathbb{E}(w_{\ell_1} \otimes M_1 \otimes w_{\ell_3}) + \mathbb{E}(M_1 \otimes w_{\ell_2} \otimes w_{\ell_3})], \\ &\quad - \frac{2c_0^3}{c_0(c_0 + 1)(c_0 + 2)} M_1 \otimes M_1 \otimes M_1 + \frac{2}{c_0(c_0 + 1)(c_0 + 2)} \sum_{k=1}^K c_k d_k \otimes d_k \otimes d_k. \end{aligned} \quad (36)$$

where \otimes denotes the tensor product.

Similarly to the GP/DICA cumulants (as discussed in Appendix B.3), moving the terms in the non-central moments (34), (35), (36), the following quantities are defined

$$(Pairs) = S := M_2 - \frac{c_0}{c_0 + 1} M_1 M_1^\top, \quad \text{LDA S-moment} \quad (37)$$

$$\begin{aligned} (Triples) = T &:= M_3 - \frac{c_0}{c_0 + 2} [\mathbb{E}(w_{\ell_1} \otimes w_{\ell_2} \otimes M_1) + \mathbb{E}(w_{\ell_1} \otimes M_1 \otimes w_{\ell_3}) + \mathbb{E}(M_1 \otimes w_{\ell_2} \otimes w_{\ell_3})] \\ &\quad + \frac{2c_0^2}{(c_0 + 1)(c_0 + 2)} M_1 \otimes M_1 \otimes M_1. \end{aligned} \quad \text{LDA T-moment} \quad (38)$$

¹²Note that another advantage of the DICA cumulants from Section 3.1 is that they do not require such a somewhat artificial condition: they are well-defined for any document length (even a document of length zero!).

¹³Note, the difference in the notation for the LDA moments in papers [1] and [3]. In [1], $M_1 = \mathbb{E}(w_{\ell_1})$, $M_2 = \mathbb{E}(w_{\ell_1} \otimes w_{\ell_2})$, and $M_3 = \mathbb{E}(w_{\ell_1} \otimes w_{\ell_2} \otimes w_{\ell_3})$. However, in [3], M_2 is equivalent to S in our notation and to *Pairs* in the notation of [1]; similarly, M_3 is T in our notation or *Triples* in the notation of [1].

Slightly abusing terminology, we refer to the entities S and T as the ‘‘LDA moments’’. They have the following diagonal structure

$$S = \frac{1}{c_0(c_0 + 1)} \sum_{k=1}^K c_k d_k d_k^\top, \quad (39)$$

$$T = \frac{2}{c_0(c_0 + 1)(c_0 + 2)} \sum_{k=1}^K c_k d_k \otimes d_k \otimes d_k. \quad (40)$$

Note however that this form of the LDA moments has a slightly different nature than the similar form (12) and (14) of the GP/DICA cumulants. Indeed, the former is the result of properties of the Dirichlet distribution, while the latter is the result of the independence of α ’s. However, one can think of the elements of a Dirichlet random vector as being almost independent (as, e.g., a Dirichlet random vector can be obtained from independent gamma variables through dividing each by their sum). Also, this closeness of the structures of the LDA moments and the GP cumulants can be explained by the closeness of the respective models as discussed in Section 2.

D.2 Asymptotically unbiased finite sample estimators for the LDA moments

Given realizations $w_{n\ell}$, $n = 1, \dots, N$, $\ell = 1, \dots, L_n$, of the token random variable w_ℓ , we now give the expressions for the finite sample estimates of S (37) and T (38) for the LDA model (and we re-write them as a function of the sample counts x_n).¹⁴ We use the notation $\widehat{\mathbb{E}}$ below to express a U-statistics empirical expectation over the token within a documents, uniformly averaged over the whole corpus. For example, $\widehat{\mathbb{E}}(w_{\ell_1} \otimes w_{\ell_2} \otimes \widehat{M}_1) := \frac{1}{N} \sum_{n=1}^N \frac{1}{L_n(L_n-1)} \sum_{\ell_1=1}^{L_n} \sum_{\substack{\ell_2=1 \\ \ell_2 \neq \ell_1}}^{L_n} w_{\ell_1} \otimes w_{\ell_2} \otimes \widehat{M}_1$.

$$\widehat{S} := \widehat{M}_2 - \frac{c_0}{c_0 + 1} \widehat{M}_1 \widehat{M}_1^\top, \quad (41)$$

$$\begin{aligned} \widehat{T} := & \widehat{M}_3 - \frac{c_0}{c_0 + 2} \left[\widehat{\mathbb{E}}(w_{\ell_1} \otimes w_{\ell_2} \otimes \widehat{M}_1) + \widehat{\mathbb{E}}(w_{\ell_1} \otimes \widehat{M}_1 \otimes w_{\ell_3}) + \widehat{\mathbb{E}}(\widehat{M}_1 \otimes w_{\ell_2} \otimes w_{\ell_3}) \right] \\ & + \frac{2c_0^2}{(c_0 + 1)(c_0 + 2)} \widehat{M}_1 \otimes \widehat{M}_1 \otimes \widehat{M}_1, \end{aligned} \quad (42)$$

where, as suggested in [3], unbiased U-statistics estimates of M_1 , M_2 and M_3 are:

$$\widehat{M}_1 := \widehat{\mathbb{E}}(w_\ell) = \frac{1}{N} \sum_{n=1}^N \frac{1}{L_n} \sum_{\ell=1}^{L_n} w_{n\ell} = \frac{1}{N} \sum_{n=1}^N [\delta_1]_n x_n = \frac{1}{N} X \delta_1, \quad (43)$$

$$\begin{aligned} \widehat{M}_2 := & \widehat{\mathbb{E}}(w_{\ell_1} w_{\ell_2}^\top) = \frac{1}{N} \sum_{n=1}^N \frac{1}{L_n(L_n-1)} \sum_{\ell_1=1}^{L_n} \sum_{\substack{\ell_2=1 \\ \ell_2 \neq \ell_1}}^{L_n} w_{n\ell_1} w_{n\ell_2}^\top \\ = & \frac{1}{N} \sum_{n=1}^N [\delta_2]_n \left(x_n x_n^\top - \sum_{\ell=1}^{L_n} w_{n\ell} w_{n\ell}^\top \right) \\ = & \frac{1}{N} \sum_{n=1}^N [\delta_2]_n (x_n x_n^\top - \text{diag}(x_n)) \\ = & \frac{1}{N} [X \text{diag}(\delta_2) X^\top - \text{diag}(X \delta_2)], \end{aligned} \quad (44)$$

$$(45)$$

¹⁴Note that because non-linear functions of \widehat{M}_1 appear in the expression for \widehat{S} (41) and \widehat{T} (42), the estimator is biased, i.e. $\mathbb{E}(\widehat{S}) \neq S$. The bias is small though: $\|\mathbb{E}(\widehat{S}) - S\| = O(1/N)$ and the estimator is asymptotically unbiased. This is in contrast with the estimator for the GP/DICA moments which is easily made unbiased.

$$\begin{aligned}
\widehat{M}_3 &:= \widehat{\mathbb{E}}(w_{\ell_1} \otimes w_{\ell_2} \otimes w_{\ell_3}) = \frac{1}{N} \sum_{n=1}^N \delta_{3n} \sum_{\ell_1=1}^{L_n} \sum_{\substack{\ell_2=1 \\ \ell_2 \neq \ell_1}}^{L_n} \sum_{\substack{\ell_3=1 \\ \ell_3 \neq \ell_1 \\ \ell_3 \neq \ell_2}}^{L_n} w_{n\ell_1} \otimes w_{n\ell_2} \otimes w_{n\ell_3} \\
&= \frac{1}{N} \sum_{n=1}^N [\delta_3]_n \left(x_n \otimes x_n \otimes x_n - \sum_{\ell=1}^{L_n} w_{n\ell} \otimes w_{n\ell} \otimes w_{n\ell} \right. \\
&\quad \left. - \sum_{\ell_1=1}^{L_n} \sum_{\substack{\ell_2=1 \\ \ell_2 \neq \ell_1}}^{L_n} (w_{n\ell_1} \otimes w_{n\ell_1} \otimes w_{n\ell_2} + w_{n\ell_1} \otimes w_{n\ell_2} \otimes w_{n\ell_1} + w_{n\ell_1} \otimes w_{n\ell_2} \otimes w_{n\ell_2}) \right) \\
&= \frac{1}{N} \sum_{n=1}^N [\delta_3]_n \left(x_n \otimes x_n \otimes x_n + 2 \sum_{m=1}^M x_{nm} (e_m \otimes e_m \otimes e_m) \right. \\
&\quad \left. - \sum_{m_1=1}^M \sum_{m_2=1}^M x_{nm_1} x_{nm_2} (e_{m_1} \otimes e_{m_1} \otimes e_{m_2} + e_{m_1} \otimes e_{m_2} \otimes e_{m_1} + e_{m_1} \otimes e_{m_2} \otimes e_{m_2}) \right). \quad (46)
\end{aligned}$$

Here, the vectors δ_1 , δ_2 and $\delta_3 \in \mathbb{R}^N$ are defined as $[\delta_1]_n := L_n^{-1}$; $[\delta_2]_n := (L_n(L_n - 1))^{-1}$, i.e. $[\delta_2]_n = \left[\binom{L_n}{2} 2! \right]^{-1}$ is the number of times to choose an ordered pair of tokens out of L_n tokens; $[\delta_3]_n := (L_n(L_n - 1)(L_n - 2))^{-1}$, i.e. $[\delta_3]_n = \left[\binom{L_n}{3} 3! \right]^{-1}$ is the number of times to choose an ordered triple of tokens out of L_n tokens. Note that the vectors δ_1 , δ_2 , and δ_3 have nothing to do with the Kronecker delta δ .

For a vector $a \in \mathbb{R}^N$, we sometimes use notation $[a]_n$ to denote its n -th element. Similarly, for a matrix $A \in \mathbb{R}^{M \times N}$ we use notation $[A]_{mn}$ to denote its (m, n) -th element.

There is a slight abuse of notation in the expressions above as w_ℓ is sometimes treated as a random variable (i.e. in $\widehat{\mathbb{E}}(w_\ell)$, $\widehat{\mathbb{E}}(w_{\ell_1} w_{\ell_2}^\top)$, etc.) and sometimes as its realization. However, the difference is clear from the context.

E Appendix. Practical aspects and implementation details

E.1 Whitening of S and dimensionality reduction

The algorithms from Section 4 require the computation of a whitening matrix W of S . Due to the similar diagonal structure ((39) and (12)) of the matrix S for both the LDA moments (37) and the GP/DICA cumulants (11), the computation of a whitening matrix is exactly the same in both cases.

By a whitening matrix, we mean a matrix $W \in \mathbb{R}^{K \times M}$ (in practice, $M \gg K$) that does not only whiten $S \in \mathbb{R}^{M \times M}$, but also reduces its dimensionality such that¹⁵ $WSW^\top = I_K$.

Let $S = U\Sigma U^\top$ be an orthogonal eigendecomposition of the symmetric matrix S . Let $\Sigma_{1:K}$ denotes the diagonal matrix that contains the largest K eigenvalues¹⁶ of S on its diagonal and let $U_{1:K}$ be a matrix with the respective eigenvalues in its columns. Then, a whitening matrix is

$$W = \Sigma_{1:K}^{\dagger 1/2} U_{1:K}^\top, \quad (47)$$

where $\Sigma_{1:K}^{\dagger 1/2}$ is a diagonal matrix constructed from $\Sigma_{1:K}$ by taking the inverse and the square root of its non-zero diagonal values († stands for the pseudo-inverse).

In practice, when only a finite sample estimator \widehat{S} of S is available, the following finite sample estimator \widehat{W} of W can be introduced

$$\widehat{W} := \widehat{\Sigma}_{1:K}^{\dagger 1/2} \widehat{U}_{1:K}^\top, \quad (48)$$

where $\widehat{S} = \widehat{U}\widehat{\Sigma}\widehat{U}^\top$.

¹⁵Note that such a whitening matrix $W \in \mathbb{R}^{K \times M}$ is not uniquely defined as left multiplication by any orthogonal matrix $V \in \mathbb{R}^{K \times K}$ does not change anything. Indeed, let $\widetilde{W} = VW$, then $\widetilde{W}S\widetilde{W}^\top = VWSW^\top V^\top = I_K$.

¹⁶We mean the largest non-negative eigenvalues. In theory, S have to be PSD. In practice, when we deal with finite number of samples, respective estimate of S can have negative eigenvalues. However, for K sufficiently small, S should have enough positive eigenvalues. Moreover, it is standard practice to use eigenvalues of S for estimation of a good value of K , e.g., by thresholding all negative and close to zero eigenvalues.

E.2 Computation of the finite sample estimators of the GP/DICA cumulants

In this section, we present efficient formulas for computation of the finite sample estimate (see Appendix B.4 for the definition of \widehat{T}) of $\widehat{W}\widehat{T}(v)\widehat{W}^\top$ for the GP/DICA models. The construction of the finite sample estimator \widehat{W} is discussed in Appendix E.1, while the computation of \widehat{S} (27) is straightforward.

By plugging the definition of the tensor \widehat{T} (28) in the formula (17) for the projection of a tensor onto a vector, we obtain for a given $v \in \mathbb{R}^M$:

$$\begin{aligned} \left[\widehat{T}(v)\right]_{m_1 m_2} &= \sum_{m_3} \widehat{\text{cum}}(x_{m_1}, x_{m_2}, x_{m_3}) v_{m_3} + 2 \sum_{m_3} \delta(m_1, m_2, m_3) \widehat{\mathbb{E}}(x_{m_3}) v_{m_3} \\ &\quad - \sum_{m_3} \delta(m_2, m_3) \widehat{\text{cov}}(x_{m_1}, x_{m_2}) v_{m_3} \\ &\quad - \sum_{m_3} \delta(m_1, m_3) \widehat{\text{cov}}(x_{m_1}, x_{m_2}) v_{m_3} \\ &\quad - \sum_{m_3} \delta(m_1, m_2) \widehat{\text{cov}}(x_{m_1}, x_{m_3}) v_{m_3} \\ &= \sum_{m_3} \widehat{\text{cum}}(x_{m_1}, x_{m_2}, x_{m_3}) v_{m_3} + 2\delta(m_1, m_2) \widehat{\mathbb{E}}(x_{m_1}) v_{m_1} \\ &\quad - \widehat{\text{cov}}(x_{m_1}, x_{m_2}) v_{m_2} - \widehat{\text{cov}}(x_{m_1}, x_{m_2}) v_{m_1} - \delta(m_1, m_2) \sum_{m_3} \widehat{\text{cov}}(x_{m_1}, x_{m_3}) v_{m_3}. \end{aligned}$$

This gives the following for the expression $\widehat{W}\widehat{T}(v)\widehat{W}^\top$:

$$\begin{aligned} \left[\widehat{W}\widehat{T}(v)\widehat{W}^\top\right]_{k_1 k_2} &= \widehat{W}_{k_1}^\top \widehat{T}(v) \widehat{W}_{k_2} \\ &= \sum_{m_1, m_2, m_3} \widehat{\text{cum}}(x_{m_1}, x_{m_2}, x_{m_3}) v_{m_3} \widehat{W}_{k_1 m_1} \widehat{W}_{k_2 m_2} \\ &\quad + 2 \sum_{m_1, m_2} \delta(m_1, m_2) \widehat{\mathbb{E}}(x_{m_1}) v_{m_1} \widehat{W}_{k_1 m_1} \widehat{W}_{k_2 m_2} \\ &\quad - \sum_{m_1, m_2} \widehat{\text{cov}}(x_{m_1}, x_{m_2}) v_{m_2} \widehat{W}_{k_1 m_1} \widehat{W}_{k_2 m_2} \\ &\quad - \sum_{m_1, m_2} \widehat{\text{cov}}(x_{m_1}, x_{m_2}) v_{m_1} \widehat{W}_{k_1 m_1} \widehat{W}_{k_2 m_2} \\ &\quad - \sum_{m_1, m_3} \widehat{\text{cov}}(x_{m_1}, x_{m_3}) v_{m_3} \widehat{W}_{k_1 m_1} \widehat{W}_{k_2 m_1}. \end{aligned}$$

where \widehat{W}_k denotes the k -th row of \widehat{W} as a column vector. By further plugging in the expressions (29) for the unbiased finite sample estimates of $\widehat{\text{cov}}$ and $\widehat{\text{cum}}$, we further get

$$\begin{aligned} \left[\widehat{W}\widehat{T}(v)\widehat{W}^\top\right]_{k_1 k_2} &= \frac{N}{(N-1)(N-2)} \sum_n \left\langle \widehat{W}_{k_1}, x_n - \widehat{\mathbb{E}}(x) \right\rangle \left\langle \widehat{W}_{k_2}, x_n - \widehat{\mathbb{E}}(x) \right\rangle \left\langle v, x_n - \widehat{\mathbb{E}}(x) \right\rangle \\ &\quad + 2 \sum_m \widehat{\mathbb{E}}(x_m) v_m \widehat{W}_{k_1 m} \widehat{W}_{k_2 m} \\ &\quad - \frac{1}{N-1} \sum_n \left\langle \widehat{W}_{k_1}, x_n - \widehat{\mathbb{E}}(x) \right\rangle \left\langle v \circ \widehat{W}_{k_2}, x_n - \widehat{\mathbb{E}}(x) \right\rangle \\ &\quad - \frac{1}{N-1} \sum_n \left\langle v \circ \widehat{W}_{k_1}, x_n - \widehat{\mathbb{E}}(x) \right\rangle \left\langle \widehat{W}_{k_2}, x_n - \widehat{\mathbb{E}}(x) \right\rangle \\ &\quad - \frac{1}{N-1} \sum_n \left\langle \widehat{W}_{k_1} \circ \widehat{W}_{k_2}, x_n - \widehat{\mathbb{E}}(x) \right\rangle \left\langle v, x_n - \widehat{\mathbb{E}}(x) \right\rangle, \end{aligned}$$

where \circ denotes the elementwise Hadamard product. Introducing the counts matrix $X \in \mathbb{R}^{M \times N}$ where each element X_{mn} is the count of the m -th word in the n -th document (note, the matrix X

contain the vector x_n in the n -th column), we further simplify the above expression

$$\begin{aligned}
\widehat{W}\widehat{T}(v)\widehat{W}^\top &= \frac{N}{(N-1)(N-2)}(\widehat{W}X)\text{diag}[X^\top v](\widehat{W}X)^\top \\
&+ \frac{N}{(N-1)(N-2)}\langle v, \widehat{\mathbb{E}}(x) \rangle \left[2N(\widehat{W}\widehat{\mathbb{E}}(x))(\widehat{W}\widehat{\mathbb{E}}(x))^\top - (\widehat{W}X)(\widehat{W}X)^\top \right] \\
&- \frac{N}{(N-1)(N-2)} \left[\widehat{W}X(X^\top v)(\widehat{W}\widehat{\mathbb{E}}(x))^\top + \widehat{W}\widehat{\mathbb{E}}(x)(\widehat{W}X(X^\top v))^\top \right] \\
&+ 2\widehat{W}\text{diag}[v \circ \widehat{\mathbb{E}}(x)]\widehat{W}^\top \\
&- \frac{1}{N-1} \left[(\widehat{W}X)(\widehat{W}\text{diag}(v)X)^\top + (\widehat{W}\text{diag}(v)X)(\widehat{W}X)^\top + \widehat{W}\text{diag}[X(X^\top v)]\widehat{W}^\top \right] \\
&+ \frac{N}{N-1} \left[(\widehat{W}\widehat{\mathbb{E}}(x))(\widehat{W}\text{diag}[v]\widehat{\mathbb{E}}(x))^\top + (\widehat{W}\text{diag}[v]\widehat{\mathbb{E}}(x))(\widehat{W}\widehat{\mathbb{E}}(x))^\top \right] \\
&+ \frac{N}{N-1} \langle v, \widehat{\mathbb{E}}(x) \rangle \widehat{W}\text{diag}[\widehat{\mathbb{E}}(x)]\widehat{W}^\top.
\end{aligned} \tag{49}$$

From expression (49), we can see that the most computationally expensive part of computing $\widehat{W}\widehat{T}(v)\widehat{W}^\top$ is the computation of the product of the whitening matrix $\widehat{W} \in \mathbb{R}^{K \times M}$ and counts matrix $X \in \mathbb{R}^{M \times N}$. As the latter is a sparse matrix, the complexity of this operation is approximately $O(R_{max}NK)$, where R_{max} is the largest number of unique words (non-zero counts) in a document. Moreover, if (49) has to be computed multiple times for different vectors $\{v_1, \dots, v_P\}$, many operations, e.g. $\widehat{W}X$, $(\widehat{W}X)(\widehat{W}X)^\top$, $(\widehat{W}\widehat{\mathbb{E}}(x))(\widehat{W}\widehat{\mathbb{E}}(x))^\top$, do not have to be recomputed P times, which makes the overall computation time significantly faster.

A more compact way to write down expression (49) is as follows

$$\begin{aligned}
\widehat{W}\widehat{T}(v)\widehat{W}^\top &= \frac{N}{(N-1)(N-2)} \left[T_1 + \langle v, \widehat{\mathbb{E}}(x) \rangle (T_2 - T_3) - (T_4 + T_4^\top) \right] \\
&+ \frac{1}{N-1} \left[T_5 + T_5^\top - T_6 - T_6^\top + \widehat{W}\text{diag}(a)\widehat{W}^\top \right],
\end{aligned} \tag{50}$$

where

$$\begin{aligned}
T_1 &= (\widehat{W}X)\text{diag}[X^\top v](\widehat{W}X)^\top, \\
T_2 &= 2N(\widehat{W}\widehat{\mathbb{E}}(x))(\widehat{W}\widehat{\mathbb{E}}(x))^\top, \\
T_3 &= (\widehat{W}X)(\widehat{W}X)^\top, \\
T_4 &= \widehat{W}X(X^\top v)(\widehat{W}\widehat{\mathbb{E}}(x))^\top, \\
T_5 &= (\widehat{W}X)(\widehat{W}\text{diag}(v)X)^\top, \\
T_6 &= (\widehat{W}\text{diag}(v)\widehat{\mathbb{E}}(x))(\widehat{W}\widehat{\mathbb{E}}(x))^\top, \\
a &= 2(N-1)[v \circ \widehat{\mathbb{E}}(x)] + \langle v, \widehat{\mathbb{E}}(x) \rangle \widehat{\mathbb{E}}(x) - X(X^\top v).
\end{aligned}$$

E.3 Computation of the finite sample estimators of the LDA moments

In this section, we present efficient formulas for computation of the finite sample estimate (see Appendix D.2 for the definition of \widehat{T}) of $\widehat{W}\widehat{T}(v)\widehat{W}^\top$ for the LDA model. Note that the construction of the sample estimator \widehat{W} of a whitening matrix W is discussed in Appendix E.1). The computation of \widehat{S} (41) is straightforward. This approach to efficient implementation was discussed in [3], however, to the best of our knowledge, the final expressions were not explicitly stated before. All derivations are straightforward, but quite tedious.

By analogy with the GP/DICA case, a projection (17) of the tensor $\widehat{T} \in \mathbb{R}^{M \times M \times M}$ (42) onto some vector $v \in \mathbb{R}^M$ in the LDA is

$$\begin{aligned}
[\widehat{T}(v)]_{m_1 m_2} &= \sum_{m_3=1}^M [\widehat{M}_3]_{m_1 m_2 m_3} v_{m_3} + \frac{2c_0^2}{(c_0+1)(c_0+2)} \sum_{m_3} [\widehat{M}_1]_{m_1} [\widehat{M}_1]_{m_2} [\widehat{M}_1]_{m_3} v_{m_3} \\
&- \frac{c_0}{c_0+2} \sum_{m_3=1}^M \left[\widehat{\mathbb{E}}(w_{\ell_1} \otimes w_{\ell_2} \otimes \widehat{M}_1) + \widehat{\mathbb{E}}(w_{\ell_1} \otimes \widehat{M}_1 \otimes w_{\ell_3}) + \widehat{\mathbb{E}}(\widehat{M}_1 \otimes w_{\ell_2} \otimes w_{\ell_3}) \right]_{m_1 m_2 m_3} v_{m_3}.
\end{aligned}$$

Plugging in the expression (46) for an unbiased sample estimate \widehat{M}_3 of M_3 , we get

$$\begin{aligned} [\widehat{T}(v)]_{m_1 m_2} &= \frac{1}{N} \sum_{n=1}^N [\delta_3]_n \left(x_{nm_1} x_{nm_2} \langle x_n, v \rangle + 2 \sum_{m_3} \delta(m_1, m_2, m_3) x_{nm_3} v_{m_3} \right) \\ &\quad - \frac{1}{N} \sum_{n=1}^N [\delta_3]_n \sum_{m_3=1}^M \left[\sum_{i,j=1}^M x_{ni} x_{nj} (e_i \otimes e_i \otimes e_j + e_i \otimes e_j \otimes e_i + e_i \otimes e_j \otimes e_j) \right]_{m_1 m_2 m_3} v_{m_3} \\ &\quad + \frac{2c_0^2}{(c_0+1)(c_0+2)} [\widehat{M}_1]_{m_1} [\widehat{M}_1]_{m_2} \langle \widehat{M}_1, v \rangle \\ &\quad - \frac{c_0}{c_0+2} \left([\widehat{M}_2]_{m_1 m_2} \langle \widehat{M}_1, v \rangle + \sum_{m_3=1}^M \left([\widehat{M}_2]_{m_1 m_3} [\widehat{M}_1]_{m_2} v_{m_3} + [\widehat{M}_2]_{m_2 m_3} [\widehat{M}_1]_{m_1} v_{m_3} \right) \right), \end{aligned}$$

where e_1, e_2, \dots, e_M denote the canonical vectors of \mathbb{R}^M (i.e. the columns of the identity matrix I_M). Further, this gives the following for the expression $\widehat{W} \widehat{T}(v) \widehat{W}^\top$:

$$\begin{aligned} [\widehat{W} \widehat{T}(v) \widehat{W}^\top]_{k_1 k_2} &= \frac{1}{N} \sum_{n=1}^N [\delta_3]_n \left(\langle x_n, v \rangle \langle x_n, \widehat{W}_{k_1} \rangle \langle x_n, \widehat{W}_{k_2} \rangle + 2 \sum_{m=1}^M x_{nm} v_m \widehat{W}_{k_1 m} \widehat{W}_{k_2 m} \right) \\ &\quad - \frac{1}{N} \sum_{n=1}^N \delta_{3n} \sum_{i,j=1}^M x_{ni} x_{nj} \left(\widehat{W}_{k_1 i} \widehat{W}_{k_2 j} v_j + \widehat{W}_{k_1 i} \widehat{W}_{k_2 j} v_i + \widehat{W}_{k_1 i} \widehat{W}_{k_2 j} v_j \right) \\ &\quad - \frac{c_0}{c_0+2} \left(\langle \widehat{W}_{k_1}, [\widehat{M}_2] \widehat{W}_{k_2} \rangle + \langle \widehat{W}_{k_1}, \widehat{M}_2 v \rangle \langle \widehat{M}_1 \widehat{W}_{k_2} \rangle + \langle \widehat{W}_{k_2}, \widehat{M}_2 v \rangle \langle \widehat{M}_1, \widehat{W}_{k_1} \rangle \right) \\ &\quad + \frac{2c_0^2}{(c_0+1)(c_0+2)} \langle \widehat{M}_1, \widehat{W}_{k_1} \rangle \langle \widehat{M}_1, \widehat{W}_{k_2} \rangle \langle \widehat{M}_1, v \rangle, \end{aligned}$$

where \widehat{W}_k denotes the k -th row of \widehat{W} as a column-vector. This further simplifies to

$$\begin{aligned} \widehat{W} \widehat{T}(v) \widehat{W}^\top &= \frac{1}{N} (\widehat{W} X) \text{diag} [(X^\top v) \circ \delta_3] (\widehat{W} X)^\top \\ &\quad + \frac{1}{N} \widehat{W} \text{diag} [2[(X \delta_3) \circ v] - X[(X^\top v) \circ \delta_3]] \widehat{W}^\top \\ &\quad - \frac{1}{N} (\widehat{W} \text{diag}[v] X) \text{diag}[\delta_3] (\widehat{W} X)^\top \\ &\quad - \frac{1}{N} (\widehat{W} X) \text{diag}[\delta_3] (\widehat{W} \text{diag}[v] X)^\top \\ &\quad - \frac{c_0}{c_0+2} \left[\langle \widehat{M}_1, v \rangle (\widehat{W} \widehat{M}_2 \widehat{W}^\top) + (\widehat{W} (\widehat{M}_2 v)) (\widehat{W} \widehat{M}_1)^\top + (\widehat{W} \widehat{M}_1) (\widehat{W} (\widehat{M}_2 v))^\top \right] \\ &\quad + \frac{2c_0^2}{(c_0+1)(c_0+2)} \langle \widehat{M}_1, v \rangle (\widehat{W} \widehat{M}_1) (\widehat{W} \widehat{M}_1)^\top. \end{aligned} \tag{51}$$

A more compact representation gives:

$$\begin{aligned} \widehat{W} \widehat{T}(v) \widehat{W}^\top &= \frac{1}{N} [T_1 + T_2 - T_3 - T_3^\top] - \frac{c_0}{c_0+2} [\langle \widehat{M}_1, v \rangle (\widehat{W} \widehat{M}_2 \widehat{W}^\top) + T_4 + T_4^\top] \\ &\quad + \frac{2c_0^2}{(c_0+1)(c_0+2)} \langle \widehat{M}_1, v \rangle (\widehat{W} \widehat{M}_1) (\widehat{W} \widehat{M}_1)^\top, \end{aligned} \tag{52}$$

where

$$\begin{aligned} T_1 &= (\widehat{W} X) \text{diag} [(X^\top v) \circ \delta_3] (\widehat{W} X)^\top, \\ T_2 &= \widehat{W} \text{diag} [2[(X \delta_3) \circ v] - X[(X^\top v) \circ \delta_3]] \widehat{W}^\top, \\ T_3 &= [\widehat{W} \text{diag}(v) X] \text{diag}(\delta_3) (\widehat{W} X)^\top, \\ T_4 &= [\widehat{W} (\widehat{M}_2 v)] (\widehat{W} \widehat{M}_1)^\top. \end{aligned}$$

The comments regarding the computational time of (50) also apply here. However, in practice we noticed that the computation of (50) is slightly faster for large datasets than the computation of (52) (although the code for both was equally well optimized). This means that the constant in $O(RNK)$ for the LDA moments is, probably, slightly larger than for the GP/DICA cumulants.

E.4 Estimation of the model parameters for GP/DICA model

Given the output A and a_p of Algorithm 1, the topic matrix is estimated as

$$\hat{d}_k := \max(0, [A^\dagger]_{:,k}) / \|\max(0, [A^\dagger]_{:,k})\|_1.$$

The truncation is necessary to enforce non-negativity and the normalization corresponds to the transformation from \tilde{D} to D , where in the latter each topic is in the simplex.

To estimate the parameters for the prior distribution over topic intensities α_k for the DICA model (5), we use the diagonalized form of the projected tensor from (18) and relate it to the output diagonal elements a_p for the p -th projection:

$$[a_p]_k = \tilde{t}_k \langle z_k, u_p \rangle = \frac{t_k}{s_k^{3/2}} \langle z_k, u_p \rangle = \frac{\text{cum}(\alpha_k, \alpha_k, \alpha_k)}{[\text{var}(\alpha_k)]^{3/2}} \langle d_k, W^\top u_p \rangle. \quad (53)$$

This formula is valid for any prior on α_k in the DICA model. For the GP model (4) where $\alpha_k \sim \text{Gamma}(c_k, b)$, we have that $\text{var}(\alpha_k) = \frac{c_k}{b^2}$ and $\text{cum}(\alpha_k, \alpha_k, \alpha_k) = \frac{2c_k}{b^3}$, and thus $\tilde{t}_k = \frac{2}{\sqrt{c_k}}$, which enables us to estimate c_k . Plugging this value of \tilde{t}_k in (53), and solving for c_k gives the following expression:

$$c_k = \frac{4 \langle d_k, W^\top u_p \rangle^2}{[a_p]_k^2}.$$

By replacing the quantities on the RHS with their estimated ones, we get one estimate for c_k per projection. We use as our final estimate the average estimate over the projections:

$$\hat{c}_k := \frac{1}{P} \sum_{p=1}^P \frac{4 \langle \hat{d}_k, \hat{W}^\top u_p \rangle^2}{[a_p]_k^2}. \quad (54)$$

Reusing the properties of the length of documents for the GP model as described in Appendix A.2, we finally use the following estimates for rate parameter b of the gamma distribution:

$$\hat{b} := \frac{\hat{c}_0}{\hat{L}}, \quad (55)$$

where $\hat{c}_0 := \sum_k \hat{c}_k$ and \hat{L} is the average document length in the corpus.

F Appendix. Supplementary experiments

F.1 The LDA moments vs parameter c_0

The LDA moments depend on the parameter c_0 , which is not trivial to set in the unsupervised setting of topic modeling, especially taking into account the complexity of evaluation for topic models [31]. In Figure 4, the joint diagonalization algorithm with the LDA moment is compared for different values of c_0 provided to the algorithm. The data is generated similarly to Figure 2. The experiment indicates that the LDA moments are somewhat sensitive to the choice of c_0 . For example, the recovery ℓ_1 -error doubles when moving from the correct choice $c_0 = 1$ to the plausible alternative $c_0 = 0.1$ for $K = 10$ on the *LDAfix1(200)* dataset (JD-LDA(10) line on the right of Figure 4). When the error is already high (for both datasets when $K = 50$ for example), then the choice of c_0 seems to matter less: the performance is uniformly bad.

F.2 Comparison of the ℓ_1 - and ℓ_2 -errors

The sample complexity results [1] for the spectral algorithm for the LDA moments allow straightforward extension to the GP/DICA cumulants, if the results from Proposition 3.1 are taken into account. The analysis is, however, in terms of the ℓ_2 -norm. Therefore, in Figure 5, we provide experimental comparison of the ℓ_1 - and ℓ_2 -errors to verify that they are indeed behaving similarly.



Figure 4: Performance of the LDA moments depending on the parameter c_0 . D and c are learned from the AP dataset for $K = 10$ and $K = 50$ and true $c_0 = 1$. JD-GP(10) for $K = 10$ and JD-GP(50) for $K = 50$. Number of sampled documents $N = 20,000$. For the error bars, each dataset is resampled 5 times. Data (left): GP sampling; (right): *LDafix1*(200) sampling. Note: a smaller value of the ℓ_1 -error is better.

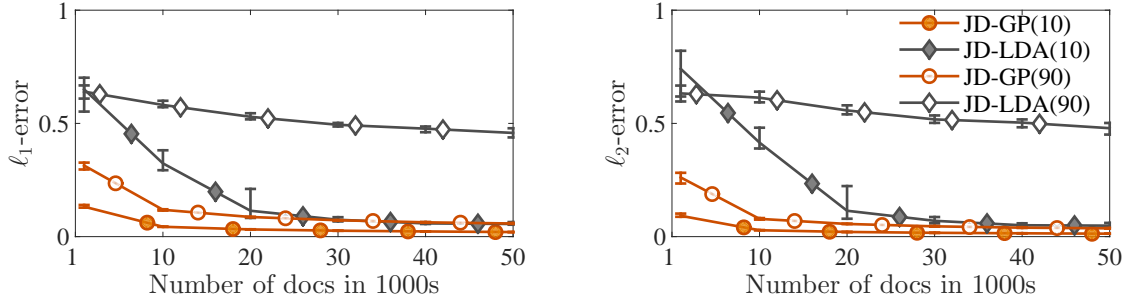


Figure 5: Comparison of the ℓ_1 - and ℓ_2 - errors on the NIPS semi-synthetic dataset as in Figure 3 (top, left). The ℓ_2 norms of the topics were normalized to 1 for the computation of the ℓ_2 error.

F.3 One more real data experiment

In Figure 6 (right), we demonstrate one more experiment with real data as described in Section 5.3. Although the variational inference outperforms the joint diagonalization algorithm, the variational inference with warm JD-restarts is the best. Note that the fact that the joint diagonalization algorithm for the LDA moments is worse than the spectral algorithm indicates that the diagonal structure (39) and (40) might not be present in the sample estimates (41) and (42).

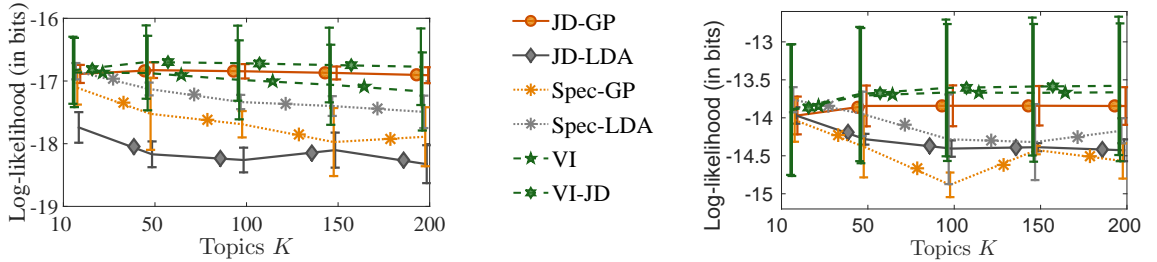


Figure 6: Experiments with real data. **Left:** the KOS dataset. **Right:** the NIPS dataset. Note: a higher value of the log-likelihood is better.

Supplementary References

- [32] A. Anandkumar, D.P. Foster, D. Hsu, S.M. Kakade, and Y.-K. Liu. A spectral algorithm for latent Dirichlet allocation. *Algorithmica*, 72(1):193–214, 2015.
- [33] B.A. Frigiyk, A. Kapila, and M.R. Gupta. Introduction to the Dirichlet distribution and related processes. Technical report, University of Washington, 2010.
- [34] T.G. Kolda and B.W. Bader. Tensor decompositions and applications. *SIAM Rev.*, 51(3):455–500, 2009.