
An Instability in Variational Inference for Topic Models

Behrooz Ghorbani¹ Hamid Javadi² Andrea Montanari^{1,3}

Abstract

Naive mean field variational methods are the state-of-the-art approach to inference in topic models. We show that these methods suffer from an instability that can produce misleading conclusions. Namely, for certain regimes of the model parameters, variational inference outputs a non-trivial decomposition into topics. However—for the same parameter values—the data contain no actual information about the true topic decomposition, and the output of the algorithm is uncorrelated with it. In particular, the estimated posterior mean is wrong, and estimated credible regions do not achieve the nominal coverage. We discuss how this instability is remedied by more accurate mean field approximations.

1. Introduction

Topic modeling (Blei, 2012) aims at extracting the latent structure from a corpus of documents (images or texts), that are represented as vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$. The key assumption is that the n documents are (approximately) convex combinations of a small number k of topics $\tilde{\mathbf{h}}_1, \dots, \tilde{\mathbf{h}}_k \in \mathbb{R}^d$. Conditional on the topics, documents are generated independently by letting $\mathbf{x}_a = (\sqrt{\beta}/d) \sum_{\ell=1}^k w_{a,\ell} \tilde{\mathbf{h}}_\ell + \mathbf{z}_a$, where the weights $\mathbf{w}_a = (w_{a,\ell})_{1 \leq \ell \leq k}$ and noise vectors \mathbf{z}_a are i.i.d. across $a \in \{1, \dots, n\}$. The coefficient $\beta \geq 0$ can be interpreted as a signal-to-noise ratio. It is also useful to introduce the matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ whose i -th row is \mathbf{x}_i , and therefore

$$\mathbf{X} = \frac{\sqrt{\beta}}{d} \mathbf{W} \mathbf{H}^\top + \mathbf{Z}, \quad (1.1)$$

where $\mathbf{W} \in \mathbb{R}^{n \times k}$ and $\mathbf{H} \in \mathbb{R}^{d \times k}$. The a -th row of \mathbf{W} , is the vector of weights \mathbf{w}_a , while the rows of \mathbf{H} will be denoted by $\mathbf{h}_i \in \mathbb{R}^k$.

¹Department of Electrical Engineering, Stanford University, CA ²Digital Signal Processing Group, Rice University, TX ³Department of Statistics, Stanford University, CA. Correspondence to: Behrooz Ghorbani <ghorbani@stanford.edu>.

Note that \mathbf{w}_a belongs to the simplex $P_1(k) = \{\mathbf{w} \in \mathbb{R}_{\geq 0}^k : \langle \mathbf{w}, \mathbf{1}_k \rangle = 1\}$. It is common to assume that its prior is Dirichlet: this class of models is known as *Latent Dirichlet Allocations*, or LDA (Blei et al., 2003). Here we will consider a symmetric Dirichlet prior, with all parameters equal to ν (which we will denote by $\text{Dir}(\nu; k)$). As for the topics \mathbf{H} , their prior distribution depends on the specific application. Here we will consider two simple examples: in the *Gaussian* case, we assume $(\tilde{\mathbf{h}}_i)_{i \leq d} \sim_{iid} \mathcal{N}(0, \mathbf{I}_k)$; in the *Dirichlet* case $(\tilde{\mathbf{h}}_i)_{i \leq d} \sim_{iid} \text{Dir}(\tilde{\nu}; k)$. Most of our discussion and explicit formulas will refer for simplicity to the Gaussian case. However, we derived analogous expressions for the Dirichlet model, and will compare with numerical simulations carried out under both distributions. Our methodology is indeed general. Finally, \mathbf{Z} will be a noise matrix with entries $(Z_{ij})_{i \in [n], j \in [d]} \sim_{iid} \mathcal{N}(0, 1/d)$.

In fully Bayesian topic models, the parameters of the Dirichlet distribution, as well as the topic distributions are themselves unknown and to be learned from data. Here we will work in an idealized setting in which they are known. We will also assume that data are in fact distributed according to the postulated generative model. Since we are studying the limitations of current approaches, our main point is only reinforced by assuming this idealized scenario.

Computing the posterior distribution of \mathbf{H}, \mathbf{W} given the data \mathbf{X} is computationally challenging. Since the seminal work of Blei, Ng and Jordan (Blei et al., 2003), variational inference is the method of choice for addressing this problem within topic models. The term ‘variational inference’ refers to a broad class of methods that aim at approximating the posterior computation by solving an optimization problem, see (Jordan et al., 1999; Wainwright et al., 2008; Blei et al., 2017) for background. A popular starting point is the Gibbs variational principle, namely the fact that the posterior solves the following convex optimization problem:

$$p_{\mathbf{W}, \mathbf{H} | \mathbf{X}}(\cdot, \cdot, | \mathbf{X}) = \arg \min_{q \in \mathcal{P}_{n,d,k}} \text{KL}(q || p_{\mathbf{W}, \mathbf{H} | \mathbf{X}}) \quad (1.2)$$

where $\text{KL}(\cdot || \cdot)$ denotes the Kullback-Leibler divergence. Optimization is within the space $\mathcal{P}_{n,d,k}$ of probability measures on \mathbf{H}, \mathbf{W} .

Even for \mathbf{W}, \mathbf{H} discrete, the Gibbs principle has exponentially many decision variables. Variational methods differ in the way the problem (1.2) is approximated. The main

approach within topic modeling is *naive mean field*, which restricts the optimization problem to the space of probability measures that factorize over the rows of \mathbf{W} , \mathbf{H} :

$$\hat{q}(\mathbf{W}, \mathbf{H}) = \prod_{i=1}^d q_i(\mathbf{h}_i) \prod_{a=1}^n \tilde{q}_a(\mathbf{w}_a). \quad (1.3)$$

By a suitable parametrization of the marginals q_i , \tilde{q}_a , this leads to an optimization problem of dimension $O((n+d)k)$, cf. Section 3. Despite being non-convex, this problem is separately convex in the $(q_i)_{i \leq d}$ and $(\tilde{q}_a)_{a \leq n}$, which naturally suggests the use of an alternating minimization algorithm which has been successfully deployed in a broad range of applications ranging from computer vision to genetics (Fei-Fei & Perona, 2005; Wang & Blei, 2011; Raj et al., 2014). We will refer to this as to the *naive mean field iteration*. Following a common use in the topics models literature, we will use the terms ‘variational inference’ and ‘naive mean field’ interchangeably.

The main result of this paper is that naive mean field presents an instability for learning Latent Dirichlet Allocations. We focus on the limit $n, d \rightarrow \infty$ with $n/d = \delta$ fixed. Hence, an LDA distribution is determined by the parameters (k, δ, ν, β) . We will show that there are regions in this parameter space such that the following two findings hold simultaneously:

No non-trivial estimator. Any estimator $\widehat{\mathbf{H}}, \widehat{\mathbf{W}}$ of the topic or weight matrices is asymptotically uncorrelated with the real model parameters \mathbf{H}, \mathbf{W} . In other words, the data do not contain enough signal to perform any strong inference.

Variational inference is randomly biased. Given the above, one would hope the Bayesian posterior to be centered on an unbiased estimate. In particular, $p(\mathbf{w}_a | \mathbf{X})$ (the posterior distribution over weights of document a) should be centered around the uniform distribution $\mathbf{w}_a = (1/k, \dots, 1/k)$. In contrast, we will show that the posterior produced by naive mean field is centered around a random distribution that is uncorrelated with the actual weights. Similarly, the posterior over topic vectors is centered around random vectors uncorrelated with the true topics.

One key argument in support of Bayesian methods is the hope that they provide a measure of uncertainty of the estimated variables. In view of this, the failure just described is particularly dangerous because it suggests some measure of certainty, although the estimates are essentially random. While the limitation of variational methods have been pointed out in the past, ours is the first case in which such an inconsistency is established rigorously in topic models.

Is there a way to eliminate this instability by using a better mean field approximation? We show that a promising approach is provided by a classical idea in statistical physics,

the Thouless-Anderson-Palmer (TAP) free energy (Thouless et al., 1977; Opper & Winther, 2001).

Variational inference via the TAP free energy. We show that the instability of naive mean field is remedied by using the TAP free energy instead of the naive mean field free energy. The latter can be optimized using an iterative scheme that is analogous to the naive mean field iteration and is known as approximate message passing (AMP).

The rest of the paper is organized as follows. Section 2 discusses a simpler example, \mathbb{Z}_2 -synchronization, which shares important features with latent Dirichlet allocations. Since calculations are fairly straightforward, this example allows to explain the main mathematical points in a simple context. We then present our main results about instability of naive mean field in Section 3, and discuss the use of TAP free energy to overcome the instability in Section 4. As mentioned above, all formal statements will refer to the Gaussian case, although we obtained analogous results for the Dirichlet case. The corresponding theoretical predictions will be compared with numerical results in the plots. Proofs will be deferred to the Supplementary Material (SM), which also contain further numerical illustrations and technical results.

1.1. Related Literature

Over the last fifteen years, topic models have been generalized to cover an impressive number of applications, including mixed membership models (Erosheva et al., 2004; Airoldi et al., 2008), dynamic topic models (Blei & Lafferty, 2006b), correlated topic models (Blei & Lafferty, 2006a; Blei et al., 2007), spatial LDA (Wang & Grimson, 2008), relational topic models (Chang & Blei, 2009), Bayesian tensor models (Zhou et al., 2015). While other approaches have been used (e.g. Gibbs sampling), variational algorithms allow to leverage advances in optimization algorithms and architectures towards the goal of variational inference (Hoffman et al., 2010; Broderick et al., 2013).

Despite this broad empirical success, little is rigorously known about the accuracy of variational inference in concrete statistical problems. Wang and Titterington (Wang & Titterington, 2004; Wang et al., 2006) studied the missing data and Gaussian mixture models. In the context of Gaussian mixtures, the same authors prove that the covariance of the variational posterior is asymptotically smaller (in the positive semidefinite order) than the inverse of the Fisher information matrix (Wang & Titterington, 2005) (see also (Giordano et al., 2015)). All of these results are established in the classical large sample asymptotics $n \rightarrow \infty$ with d fixed. In the present paper we focus instead on the high-dimensional limit $n = \Theta(d)$ and prove that also the mode (or mean) of the variational posterior is incorrect. The high-dimensional regime is particularly relevant for

the analysis of Bayesian methods. Indeed, in the classical low-dimensional asymptotics Bayesian approaches do not outperform maximum likelihood.

Naive mean field variational inference was used in (Celisse et al., 2012; Bickel et al., 2013) to estimate the parameters of the stochastic block model. The recent paper (Zhang & Zhou, 2017) also studies variational inference in the context of the stochastic block model. The work of (Celisse et al., 2012; Bickel et al., 2013; Zhang & Zhou, 2017) establishes positive results at large signal-to-noise ratio (albeit for a different model), while we prove inconsistency at low signal-to-noise ratio. General conditions for consistency of variational Bayes methods are proposed in (Pati et al., 2017).

Our work also builds on recent theoretical advances in high-dimensional low-rank models, that were mainly driven by techniques from mathematical statistical physics (more specifically, spin glass theory). An incomplete list of relevant references includes (Korada & Macris, 2009; Deshpande & Montanari, 2014; Deshpande et al., 2016; Krzakala et al., 2016; Barbier et al., 2016; Lelarge & Miolane, 2016; Miolane, 2017; Lesieur et al., 2017; Alaoui & Krzakala, 2018). These papers prove asymptotically exact characterizations of the Bayes optimal estimation error in low-rank models, to an increasing degree of generality, under the high-dimensional scaling $n, d \rightarrow \infty$ with $n/d \rightarrow \delta \in (0, \infty)$.

Related ideas also suggest an iterative algorithm for Bayesian estimation, namely Bayes Approximate Message Passing (Donoho et al., 2009; 2010). As mentioned above, Bayes AMP can be regarded as minimizing a different variational approximation known as the TAP free energy. An important advantage over naive mean field is that AMP can be rigorously analyzed using a method known as state evolution (Bayati & Montanari, 2011; Javanmard & Montanari, 2013; Berthier et al., 2017).

1.2. Notations

We denote by I_m the identity matrix, and by J_m the all-ones matrix in m dimensions. We use $\mathbf{1}_k \in \mathbb{R}^k$ for the all-ones vector. We will use \otimes for the tensor (outer) product. In particular, given vectors expressed in the canonical basis as $\mathbf{u} = \sum_{i=1}^{d_1} u_i \mathbf{e}_i \in \mathbb{R}^{d_1}$ and $\mathbf{v} = \sum_{j=1}^{d_2} v_j \mathbf{e}_j \in \mathbb{R}^{d_2}$, $\mathbf{u} \otimes \mathbf{v} \in \mathbb{R}^{d_1} \otimes \mathbb{R}^{d_2}$ is the tensor with coordinates $(\mathbf{u} \otimes \mathbf{v})_{ij} = u_i v_j$ in the basis $\mathbf{e}_i \otimes \mathbf{e}_j$. We will identify the space of matrices $\mathbb{R}^{d_1 \times d_2}$ with the tensor product $\mathbb{R}^{d_1} \otimes \mathbb{R}^{d_2}$. Given a symmetric matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$, we denote by $\lambda_1(\mathbf{M}) \geq \lambda_2(\mathbf{M}) \geq \dots \geq \lambda_n(\mathbf{M})$ its eigenvalues in decreasing order. For a matrix (or vector) $\mathbf{A} \in \mathbb{R}^{d \times n}$ we denote the orthogonal projector onto the subspace spanned by the columns of \mathbf{A} by $\mathbf{P}_A \in \mathbb{R}^{d \times d}$, and its orthogonal complement by $\mathbf{P}_A^\perp = I_d - \mathbf{P}_A$. When the subscript is

omitted, it is understood that $\mathbf{P} = \mathbf{1}_d \mathbf{1}_d / d$ and $\mathbf{P}_\perp = I_d - \mathbf{P}$.

2. A Toy Example: \mathbb{Z}_2 -Synchronization

Before passing to the main results, it is useful to present the main ideas on a toy example. In \mathbb{Z}_2 synchronization we are interested in estimating a vector $\boldsymbol{\sigma} \in \{+1, -1\}^n$ from observations $\mathbf{X} \in \mathbb{R}^{n \times n}$, generated according to

$$\mathbf{X} = \frac{\lambda}{n} \boldsymbol{\sigma} \boldsymbol{\sigma}^\top + \mathbf{Z}, \quad (2.1)$$

where $\mathbf{Z} = \mathbf{Z}^\top \in \mathbb{R}^{n \times n}$ is a noise matrix from the Gaussian Orthogonal Ensemble $\text{GOE}(n)$, namely $(Z_{ij})_{i < j \leq n} \sim_{iid} \mathcal{N}(0, 1/n)$ are independent of $(Z_{ii})_{i \leq n} \sim_{iid} \mathcal{N}(0, 2/n)$. The parameter $\lambda \geq 0$ corresponds to the signal-to-noise ratio.

It is known that for $\lambda \leq 1$ no algorithm can estimate $\boldsymbol{\sigma}$ from data \mathbf{X} with positive correlation in the limit $n \rightarrow \infty$. The following is an immediate consequence of (Korada & Macris, 2009; Deshpande et al., 2016), see Supplementary Material (SM).

Lemma 2.1. *Under model (2.1), for $\lambda \leq 1$ and any estimator $\hat{\boldsymbol{\sigma}} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^n \setminus \{\mathbf{0}\}$, the following limit holds in probability:*

$$\limsup_{n \rightarrow \infty} \frac{|\langle \hat{\boldsymbol{\sigma}}(\mathbf{X}), \boldsymbol{\sigma} \rangle|}{\|\hat{\boldsymbol{\sigma}}(\mathbf{X})\|_2 \|\boldsymbol{\sigma}\|_2} = 0. \quad (2.2)$$

How does variational inference perform on this problem? Any product probability distribution $\hat{q}(\boldsymbol{\sigma}) = \prod_{i=1}^n q_i(\sigma_i)$ can be parametrized by the means $m_i = \sum_{\sigma_i \in \{+1, -1\}} q_i(\sigma_i) \sigma_i$, and it is immediate to get $\text{KL}(\hat{q} \| p_{\boldsymbol{\sigma} | \mathbf{X}}) = \mathcal{F}(\mathbf{m}) + \text{const.}$, where

$$\mathcal{F}(\mathbf{m}) \equiv -\frac{\lambda}{2} \langle \mathbf{m}, \mathbf{X}_0 \mathbf{m} \rangle - \sum_{i=1}^n h(m_i). \quad (2.3)$$

Here \mathbf{X}_0 is obtained from \mathbf{X} by setting the diagonal entries to 0, and $h(x) = -\frac{(1+x)}{2} \log \frac{(1+x)}{2} - \frac{(1-x)}{2} \log \frac{(1-x)}{2}$ is the binary entropy function. In view of Lemma 2.1, the correct posterior distribution should be essentially uniform, resulting in \mathbf{m} vanishing. Indeed, $\mathbf{m}_* = \mathbf{0}$ is a stationary point of the mean field free energy $\mathcal{F}(\mathbf{m})$: $\nabla \mathcal{F}(\mathbf{m})|_{\mathbf{m}=\mathbf{m}_*} = \mathbf{0}$. We refer to this as the ‘uninformative fixed point’.

Is \mathbf{m}_ a local minimum?* Computing the Hessian at the uninformative fixed point yields

$$\nabla^2 \mathcal{F}(\mathbf{m})|_{\mathbf{m}=\mathbf{m}_*} = -\lambda \mathbf{X}_0 + \mathbf{I}. \quad (2.4)$$

The matrix \mathbf{X}_0 is a rank-one deformation of a Wigner matrix and its spectrum is well understood (Baik et al., 2005; Féral & Péché, 2007; Benaych-Georges & Nadakuditi, 2011). For

$\lambda \leq 1$, its eigenvalues are contained with high probability in the interval $[-2, 2]$, with $\lambda_{\min}(\mathbf{X}) \rightarrow -2$, $\lambda_{\max}(\mathbf{X}) \rightarrow 2$ as $n \rightarrow \infty$. For $\lambda > 1$, $\lambda_{\max}(\mathbf{X}) \rightarrow \lambda + \lambda^{-1}$, while the other eigenvalues are contained in $[-2, 2]$. This implies

$$\lim_{n \rightarrow \infty} \lambda_{\min}(\nabla^2 \mathcal{F}|_{\mathbf{m}_*}) = \begin{cases} 1 - 2\lambda & \text{if } \lambda \leq 1, \\ -\lambda^2 & \text{if } \lambda > 1. \end{cases} \quad (2.5)$$

In other words, $\mathbf{m}_* = 0$ is a local minimum for $\lambda < 1/2$, but becomes a saddle point for $\lambda > 1/2$. In particular, for $\lambda \in (1/2, 1)$, variational inference will produce an estimate $\hat{\mathbf{m}} \neq 0$, although the posterior should be essentially uniform. In fact, it is possible to make this conclusion more quantitative. In the Supplementary Material we prove that any local minimum $\hat{\mathbf{m}}$ has norm $\|\hat{\mathbf{m}}\|_2^2 \geq c_0 n$, with high probability.

The above mathematical phenomenon implies that naive mean field leads to incorrect inferential statements for $\lambda \in (1/2, 1)$. In order to formalize this point, given any estimators $\{\hat{q}_i(\cdot)\}_{i \leq n}$ of the posterior marginals, we define the per-coordinate expected coverage as

$$\mathcal{Q}(\hat{q}) = \frac{1}{n} \sum_{i=1}^n \mathbb{P}(\sigma_i = \arg \max_{\tau_i \in \{+1, -1\}} \hat{q}_i(\tau_i)). \quad (2.6)$$

This is the expected fraction of coordinates that are estimated correctly by choosing σ according to the estimated posterior. On the other hand, if the \hat{q}_i were accurate, Bayesian theory would suggest claiming the coverage

$$\hat{\mathcal{Q}}(\hat{q}) \equiv \frac{1}{n} \sum_{i \leq n} \max_{\tau_i} \hat{q}_i(\tau_i). \quad (2.7)$$

The following result shows that mean field overestimates the coverage achieved.

Theorem 1. *Let $\hat{\mathbf{m}} \in [-1, 1]^n$ be any local minimum of the mean field free energy $\mathcal{F}(\mathbf{m})$, under the \mathbb{Z}_2 -synchronization model (2.1), and consider the corresponding posterior marginal estimates $\hat{q}_i(\sigma_i) = (1 + \hat{m}_i \sigma_i)/2$. Then, there exists a numerical constant $c_0 > 0$ such that, with high probability, for $\lambda \in (1/2, 1)$,*

$$\mathcal{Q}(\hat{q}) \leq \frac{1}{2} + o_n(1), \quad \hat{\mathcal{Q}}(\hat{q}) \geq \frac{1}{2} + c_0 \min((2\lambda - 1), 1).$$

3. Instability for Topic Models

3.1. Information-Theoretic Limit

As in the case of \mathbb{Z}_2 synchronization discussed in Section 2, we expect it to be impossible to estimate the factors \mathbf{W} , \mathbf{H} with strictly positive correlation for small enough signal-to-noise ratio β (or small enough sample size δ). The exact threshold was characterized recently in (Miolane, 2017) (but see also (Deshpande & Montanari, 2014; Barbier et al.,

2016; Lelarge & Miolane, 2016; Lesieur et al., 2017) for closely related results). The characterization in (Miolane, 2017) is given in terms of a variational principle over $k \times k$ matrices.

Theorem 2 (Special case of (Miolane, 2017)). *Let $I_n(\mathbf{X}; \mathbf{W}, \mathbf{H})$ denote the mutual information between the data \mathbf{X} and the factors \mathbf{H}, \mathbf{W} under the LDA model (1.1). Then, the following limit holds almost surely*

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} I_n(\mathbf{X}; \mathbf{W}, \mathbf{H}) = \inf_{\mathbf{M} \in \mathbb{S}_k} \text{RS}(\mathbf{M}; k, \delta, \nu), \quad (3.1)$$

where \mathbb{S}_k is the cone of $k \times k$ positive semidefinite matrices and $\text{RS}(\dots)$ is a function given explicitly in SM.

It is also shown in SM that $\mathbf{M}^* = (\delta\beta/k^2)\mathbf{J}_k$ is a stationary point of the free energy $\text{RS}(\mathbf{M}; k, \delta, \nu)$. We shall refer to \mathbf{M}^* as the uninformative point. Let $\beta_{\text{Bayes}} = \beta_{\text{Bayes}}(k, \delta, \nu)$ be the supremum value of β such that the infimum in Eq. (3.1) is uniquely achieved at \mathbf{M}^* (namely, the supremum β such that $\text{RS}(\mathbf{M}; k, \delta, \nu) > \text{RS}(\mathbf{M}_*; k, \delta, \nu)$ for all $\mathbf{M} \neq \mathbf{M}_*$).

As formalized below, for $\beta < \beta_{\text{Bayes}}$ the data \mathbf{X} do not contain sufficient information for estimating \mathbf{H}, \mathbf{W} in a non-trivial manner.

Proposition 3.1. *Let $\mathbf{M}_* = \delta\beta\mathbf{J}_k/k^2$. Then \mathbf{M}^* is a stationary point of the function $\mathbf{M} \mapsto \text{RS}(\mathbf{M}; \beta, k, \delta, \nu)$. Further, it is a local minimum provided $\beta < \beta_{\text{spect}}(k, \delta, \nu)$ where the spectral threshold is given by $\beta_{\text{spect}} \equiv k(k\nu + 1)/\sqrt{\delta}$.*

Finally, if $\beta < \beta_{\text{Bayes}}(k, \delta, \nu)$, there is no estimator $\mathbf{X} \mapsto \hat{\mathbf{F}}_n(\mathbf{X})$ whose mean square error $\mathbb{E}\{\|\mathbf{W}\mathbf{H}^\top - \hat{\mathbf{F}}_n(\mathbf{X})\|_F^2\}$ is asymptotically smaller than the mean square error of the trivial estimator $\hat{\mathbf{F}}_n(\mathbf{X}) = c\mathbf{1}_n(\mathbf{X}^\top \mathbf{1}_n)^\top$, for $c \equiv \sqrt{\beta}/(k + \beta\delta)$ a constant.

This result compares the mean square error of an arbitrary estimator $\hat{\mathbf{F}}_n$, to the mean square error of the trivial estimator that replaces each column of \mathbf{X} by its average. Of course, $\beta_{\text{Bayes}} \leq \beta_{\text{spect}}$. However, this upper bound appears to be tight for small k .

Remark 3.1. Solving numerically the $k(k + 1)/2$ -dimensional problem (3.1) indicates that $\beta_{\text{Bayes}}(k, \nu, \delta) = \beta_{\text{spect}}(k, \nu, \delta)$ for $k \in \{2, 3\}$ and $\nu = 1$.

3.2. Naive Mean Field Free Energy

We consider a trial joint distribution that factorizes according to rows of \mathbf{W} and \mathbf{H} according to Eq. (1.3). It turns out (see SM) that, for any stationary point of $\text{KL}(\hat{q} \| p_{\mathbf{H}, \mathbf{W} | \mathbf{X}})$ over such product distributions, the marginals take the form

$$q_i(\mathbf{h}) = e^{\langle \mathbf{m}_i, \mathbf{h} \rangle - \frac{1}{2} \langle \mathbf{h}, \mathbf{Q}_i \mathbf{h} \rangle - \phi(\mathbf{m}_i, \mathbf{Q}_i)} q_0(\mathbf{h}), \quad (3.2)$$

where $q_0(\cdot)$ is the prior distribution of \mathbf{h}_i (the i -th row of \mathbf{H}), and $\phi : \mathbb{R}^k \times \mathbb{R}^{k \times k} \rightarrow \mathbb{R}$ is defined implicitly by the normalization condition $\int q_i(d\mathbf{h}_i) = 1$. A similar form holds for $\tilde{q}_a(\mathbf{w})$, with parameters $\tilde{\mathbf{m}}_a, \tilde{\mathbf{Q}}_a$, and normalization factor $\tilde{\phi}(\tilde{\mathbf{m}}_a, \tilde{\mathbf{Q}}_a)$. In the following we let $\mathbf{m} = (\mathbf{m}_i)_{i \leq d}, \tilde{\mathbf{m}} = (\tilde{\mathbf{m}}_a)_{a \leq n}$ denote the set of parameters in these distributions; these can also be viewed as matrices $\mathbf{m} \in \mathbb{R}^{d \times k}$ and $\tilde{\mathbf{m}} \in \mathbb{R}^{d \times k}$ whose i -th row is \mathbf{m}_i (in the former case) or $\tilde{\mathbf{m}}_i$ (in the latter).

It is useful to define the functions $F, \tilde{F} : \mathbb{R}^k \times \mathbb{R}^{k \times k} \rightarrow \mathbb{R}^k$ and $G, \tilde{G} : \mathbb{R}^k \times \mathbb{R}^{k \times k} \rightarrow \mathbb{R}^{k \times k}$ as (proportional to) expectations with respect to the approximate posteriors (3.2)

$$F(\mathbf{m}_i; \mathbf{Q}) \equiv \sqrt{\beta} \int \mathbf{h} q_i(d\mathbf{h}), \quad (3.3)$$

$$G(\mathbf{m}_i; \mathbf{Q}) \equiv \beta \int \mathbf{h} \otimes \mathbf{h} q_i(d\mathbf{h}). \quad (3.4)$$

Similarly $\tilde{F}(\tilde{\mathbf{m}}_a; \tilde{\mathbf{Q}}), \tilde{G}(\tilde{\mathbf{m}}_a; \tilde{\mathbf{Q}})$ will denote the first and second moments of $\tilde{q}_a(\mathbf{w})$. For $\mathbf{m} \in \mathbb{R}^{d \times k}$, we overload the notation and denote by $F(\mathbf{m}; \mathbf{Q}) \in \mathbb{R}^{d \times k}$ the matrix whose i -th row is $F(\mathbf{m}_i; \mathbf{Q})$ (and similarly for $\tilde{F}(\tilde{\mathbf{m}}; \tilde{\mathbf{Q}})$).

When restricted to a product-form ansatz with parametrization (3.2), the mean field free energy takes the form (see SM) $\text{KL}(\hat{q} \| p_{\mathbf{W}, \mathbf{H} | \mathbf{X}}) = \mathcal{F}(\mathbf{r}, \tilde{\mathbf{r}}, \Omega, \tilde{\Omega}) + d \|\mathbf{X}\|_F^2 / 2 + \log p_{\mathbf{X}}(\mathbf{X})$, where

$$\begin{aligned} \mathcal{F}(\mathbf{r}, \tilde{\mathbf{r}}, \Omega, \tilde{\Omega}) &= \sum_{i=1}^d \psi_*(\mathbf{r}_i, \Omega_i) + \sum_{a=1}^n \tilde{\psi}_*(\tilde{\mathbf{r}}_a, \tilde{\Omega}_a) \\ &- \sqrt{\beta} \text{Tr}(\mathbf{X} \mathbf{r} \tilde{\mathbf{r}}^\top) + \frac{\beta}{2d} \sum_{i=1}^d \sum_{a=1}^n \langle \Omega_i, \tilde{\Omega}_a \rangle, \end{aligned} \quad (3.5)$$

and $\psi_*, \tilde{\psi}_*$ are the Legendre duals of $\phi, \tilde{\phi}$, e.g.

$$\psi_*(\mathbf{r}, \Omega) \equiv \sup_{\mathbf{m}, \mathbf{Q}} \left\{ \langle \mathbf{r}, \mathbf{m} \rangle - \frac{1}{2} \langle \Omega, \mathbf{Q} \rangle - \phi(\mathbf{m}, \mathbf{Q}) \right\}.$$

This equation implies a convex duality relation between $(\mathbf{r}, \tilde{\mathbf{r}}, \Omega, \tilde{\Omega})$ and $(\mathbf{m}, \tilde{\mathbf{m}}, \mathbf{Q}, \tilde{\mathbf{Q}})$. Namely

$$\mathbf{r}_i \equiv \frac{1}{\sqrt{\beta}} F(\mathbf{m}_i; \mathbf{Q}), \quad \Omega_i \equiv \frac{1}{\beta} G(\mathbf{m}_i; \mathbf{Q}), \quad (3.6)$$

and similarly for $\tilde{\mathbf{r}}_a, \tilde{\Omega}_a$ and $\tilde{\mathbf{m}}_a, \tilde{\mathbf{Q}}_a$. By strict convexity of $\phi(\mathbf{m}, \mathbf{Q}), \tilde{\phi}(\tilde{\mathbf{m}}, \tilde{\mathbf{Q}})$ (the latter is strongly convex on the hyperplane $\langle \mathbf{1}, \tilde{\mathbf{m}} \rangle = 0, \langle \mathbf{1}, \tilde{\mathbf{Q}} \mathbf{1} \rangle = 0$) we can view $\mathcal{F}(\dots)$ as a function of $(\mathbf{r}, \tilde{\mathbf{r}}, \Omega, \tilde{\Omega})$ or $(\mathbf{m}, \tilde{\mathbf{m}}, \mathbf{Q}, \tilde{\mathbf{Q}})$. With an abuse of notation, we will write $\mathcal{F}(\mathbf{r}, \tilde{\mathbf{r}}, \Omega, \tilde{\Omega})$ or $\mathcal{F}(\mathbf{m}, \tilde{\mathbf{m}}, \mathbf{Q}, \tilde{\mathbf{Q}})$ interchangeably.

A critical (stationary) point of the free energy (3.5) is a point at which $\nabla \mathcal{F}(\mathbf{m}, \tilde{\mathbf{m}}, \mathbf{Q}, \tilde{\mathbf{Q}}) = \mathbf{0}$. It turns out that the mean field free energy always admits a point that does not

distinguish between the k latent factors, and in particular $\mathbf{m} = \mathbf{v} \mathbf{1}_k^\top, \tilde{\mathbf{m}} = \tilde{\mathbf{v}} \mathbf{1}_k^\top$, as stated in detail below. We will refer to this as the *uninformative critical point* (or *uninformative fixed point*).

Lemma 3.2. *The naive mean field free energy of Eq. (3.5) admits a stationary point whereby, for all $i \in [d], a \in [n]$,*

$$\mathbf{m}_i^* = \frac{\sqrt{\beta}}{k} (\mathbf{X}^\top \mathbf{1}_n)_i \mathbf{1}_k, \quad (3.7)$$

$$\tilde{\mathbf{m}}_a^* = \frac{\beta}{k(1 + q_1^* + kq_2^*)} (\mathbf{X} \mathbf{X}^\top \mathbf{1}_n)_a \mathbf{1}_k, \quad (3.8)$$

and further $\mathbf{Q}_i^* = q_1^* \mathbf{I}_k + q_2^* \mathbf{J}_k, \tilde{\mathbf{Q}}_a^* = \tilde{q}_1^* \mathbf{I}_k + \tilde{q}_2^* \mathbf{J}_k$. The parameters q_i^*, \tilde{q}_i^* are explicitly given in the Supplementary Material.

We note that there appear always to be a unique stationary point of the form given by this lemma. Although we do not have a proof of uniqueness, in the SM we prove that the solution is unique conditional on a certain inequality that can be easily checked numerically.

3.3. Naive Mean Field Iteration

As mentioned in the introduction, the variational approximation of the free energy is often minimized by alternating minimization over the marginals $(q_i)_{i \leq d}, (\tilde{q}_a)_{a \leq n}$ of Eq. (1.3). Using the parametrization (3.2), we obtain the following naive mean field iteration for $\mathbf{m}^t, \tilde{\mathbf{m}}^t, \mathbf{Q}^t, \tilde{\mathbf{Q}}^t$ (see SM):

$$\begin{aligned} \mathbf{m}^{t+1} &= \mathbf{X}^\top \tilde{F}(\tilde{\mathbf{m}}^t; \tilde{\mathbf{Q}}^t), \quad \mathbf{Q}^{t+1} = \frac{1}{d} \sum_{a=1}^n \tilde{G}(\tilde{\mathbf{m}}_a^t; \tilde{\mathbf{Q}}^t), \\ \tilde{\mathbf{m}}^t &= \mathbf{X} F(\mathbf{m}^t; \mathbf{Q}^t), \quad \tilde{\mathbf{Q}}^t = \frac{1}{d} \sum_{i=1}^d G(\mathbf{m}_i^t; \mathbf{Q}^t). \end{aligned}$$

Note that, while the free energy naturally depends on the $(\mathbf{Q}_i)_{i \leq d}, (\tilde{\mathbf{Q}}_a)_{a \leq n}$, the iteration sets $\mathbf{Q}_i^t = \mathbf{Q}^t, \tilde{\mathbf{Q}}_a^t = \tilde{\mathbf{Q}}^t$, independent of the indices i, a . In fact, any stationary point of $\mathcal{F}(\mathbf{m}, \tilde{\mathbf{m}}, \mathbf{Q}, \tilde{\mathbf{Q}})$ can be shown to be of this form.

The state of the naive mean field iteration is given by the pair $(\mathbf{m}^t, \mathbf{Q}^t) \in \mathbb{R}^{d \times k} \times \mathbb{R}^{k \times k}$, and $(\tilde{\mathbf{m}}^t, \tilde{\mathbf{Q}}^t)$ can be viewed as derived variables. The iteration hence defines a mapping $\mathcal{M}_{\mathbf{X}} : \mathbb{R}^{d \times k} \times \mathbb{R}^{k \times k} \rightarrow \mathbb{R}^{d \times k} \times \mathbb{R}^{k \times k}$, and we can write it in the form $(\mathbf{m}^{t+1}, \mathbf{Q}^{t+1}) = \mathcal{M}_{\mathbf{X}}(\mathbf{m}^t, \mathbf{Q}^t)$. Any critical point of the free energy (3.5) is a fixed point of the naive mean field iteration and vice-versa, as shown in the SM. In particular, the uninformative critical point $(\mathbf{m}^*, \tilde{\mathbf{m}}^*, \mathbf{Q}^*, \tilde{\mathbf{Q}}^*)$ is a fixed point of the naive mean field iteration.

3.4. Instability

In view of Section 3.1, for $\beta < \beta_{\text{Bayes}}(k, \delta, \nu)$, the real posterior should be centered around a point symmetric under

permutations of the topics. In particular, the posterior $\tilde{q}(\mathbf{w}_a)$ over the weights of document a should be centered around the symmetric distribution $\mathbf{w}_a = (1/k, \dots, 1/k)$.

A minimum consistency condition for variational inference is that the uninformative stationary point is a local minimum of the posterior for $\beta < \beta_{\text{Bayes}}$. The next theorem provides a necessary condition for stability of the uninformative point, which we expect to be tight. As discussed below, it implies that this point is a saddle in an interval of β below β_{Bayes} . We recall that the index of a smooth function f at stationary point \mathbf{x}_* is the number of the negative eigenvalues of the Hessian $\nabla^2 f(\mathbf{x}_*)$.

Theorem 3. Define q_1^* , q_2^* as in Lemma 3.2, and let

$$L(\beta, k, \delta, \nu) \equiv \frac{\beta(1 + \sqrt{\delta})^2}{1 + q_1^*} \left(\frac{q_1^*}{\delta\beta} + k \left[\frac{q_2^*}{1 + q_1^* + kq_2^*} \left(\frac{1}{\delta\beta} + \frac{1}{k} \right) - \frac{1}{k^2} \right]_+ \right).$$

If $L(\beta, k, \delta, \nu) > 1$, then there exists $\varepsilon_1, \varepsilon_2 > 0$ such that the uninformative critical point of Lemma 3.2, $(\mathbf{m}^*, \tilde{\mathbf{m}}^*, \mathbf{Q}^*, \tilde{\mathbf{Q}}^*)$ is, with high probability, a saddle point, with index at least $n\varepsilon_1$ and $\lambda_{\min}(\mathcal{F}|_{\mathbf{m}^*, \tilde{\mathbf{m}}^*, \mathbf{Q}^*, \tilde{\mathbf{Q}}^*}) \leq -\varepsilon_2$.

Correspondingly $(\mathbf{m}^*, \mathbf{Q}^*)$ is an unstable critical point of the mapping $\mathcal{M}_{\mathbf{X}}$ in the sense that the Jacobian $D\mathcal{M}_{\mathbf{X}}$ has spectral radius larger than one at $(\mathbf{m}^*, \mathbf{Q}^*)$.

Remark 3.2. We established an analogous instability phase transition for the Dirichlet case. The corresponding prediction is reported in Figure 3. Explicit formulas are reported in the SM.

In the following, we will say that a fixed point $(\mathbf{m}^*, \mathbf{Q}^*)$ is stable if the linearization of $\mathcal{M}_{\mathbf{X}}(\cdot)$ at $(\mathbf{m}^*, \mathbf{Q}^*)$ (i.e. the Jacobian matrix $D\mathcal{M}_{\mathbf{X}}(\mathbf{m}^*, \mathbf{Q}^*)$) has spectral radius smaller than one. By the Hartman-Grobman linearization theorem (Perko, 2013), this implies that $(\mathbf{m}^*, \mathbf{Q}^*)$ is an attractive fixed point. Vice-versa, we say that $(\mathbf{m}^*, \mathbf{Q}^*)$ is unstable if the Jacobian $D\mathcal{M}_{\mathbf{X}}(\mathbf{m}^*, \mathbf{Q}^*)$ has spectral radius larger than one. In this case, for any neighborhood of $(\mathbf{m}^*, \mathbf{Q}^*)$, and a generic initialization in that neighborhood, $(\mathbf{m}^t, \mathbf{Q}^t)$ does not converge to the fixed point.

Motivated by Theorem 3, we define the instability threshold $\beta_{\text{inst}} = \beta_{\text{inst}}(k, \delta, \nu)$ as the infimum $\beta \geq 0$ such that $L(\beta, k, \delta, \nu) > 1$. Let us emphasize that, while we discuss the consequences of the instability at β_{inst} on the naive mean field iteration, this is a problem of the variational free energy (3.5) and not of the specific optimization algorithm.

3.5. Numerical Results for Naive Mean Field

In order to investigate the impact of the instability described above, we carried out extensive numerical simulations with

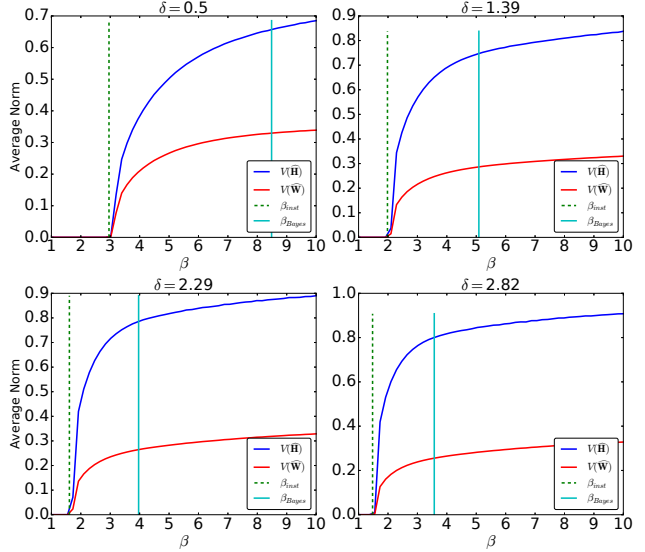


Figure 1. Normalized distances $V(\widehat{\mathbf{H}})$, $V(\widehat{\mathbf{W}})$ of the naive mean field estimates from the uninformative fixed point. Here $k = 2$, $d = 1000$ and $n = d\delta$: each data point corresponds to an average over 400 random realizations.

the naive mean field iteration. After any number of iterations t , estimates of the factors \mathbf{H} , \mathbf{W} are obtained by computing expectations with respect to the marginals (3.2). This results in

$$\widehat{\mathbf{H}}^t = \mathbf{r}^t = \frac{1}{\sqrt{\beta}} \mathbf{F}(\mathbf{m}^t; \mathbf{Q}_t), \quad \widehat{\mathbf{W}}^t = \tilde{\mathbf{r}}^t = \frac{1}{\sqrt{\beta}} \tilde{\mathbf{F}}(\tilde{\mathbf{m}}^t; \tilde{\mathbf{Q}}_t).$$

We select a two-dimensional grid of (δ, β) 's and generate 400 different instances according to the LDA model for each grid point. We report various statistics of the estimates aggregated over the 400 instances. We have performed the simulations for $\nu = \tilde{\nu} = 1$ and $k \in \{2, 3\}$, both for the Gaussian and the Dirichlet models. For space considerations, we focus here on the case $\nu = 1$, $k = 2$, and discuss other results in the SM. (Simulations for other values of ν also yield similar results.)

We initialize the naive mean field iteration near the uninformative fixed-point and iterate until a convergence criterion or the maximum number of 300 iterations is reached.

Recall the definition $\mathbf{P}_{\perp} = \mathbf{I}_k - \mathbf{1}_k \mathbf{1}_k^T / k$. In order to investigate the instability of Theorem 3, we define the quantities

$$V(\widehat{\mathbf{W}}) \equiv \frac{1}{\sqrt{n}} \|\widehat{\mathbf{W}} \mathbf{P}_{\perp}\|_F, \quad V(\widehat{\mathbf{H}}) \equiv \frac{1}{\sqrt{d}} \|\widehat{\mathbf{H}} \mathbf{P}_{\perp}\|_F$$

In Figure 1 we plot empirical results for the average $V(\widehat{\mathbf{W}})$, $V(\widehat{\mathbf{H}})$ for $k = 2$, $\nu = 1$ and four values of δ , within the Gaussian model. In Figure 2 (left frame), we plot the empirical probability that variational inference does not converge to the uninformative fixed point or, more precisely,

$\widehat{\mathbb{P}}(V(\widehat{\mathbf{W}}) \geq \varepsilon_0)$ with $\varepsilon_0 = 10^{-4}$, evaluated on a grid of (β, δ) values, for the same model. We also plot the Bayes threshold β_{Bayes} (which numerically coincides with the spectral threshold β_{spect}) and the instability β_{inst} .

It is clear from Figures 1, 2 (left frame), that variational inference stops converging to the uninformative fixed point (although we initialize close to it) when β is still significantly smaller than the Bayes threshold β_{Bayes} (i.e. in a regime in which the uninformative fixed point would be a reasonable output). The data are consistent with the hypothesis that variational inference becomes unstable at β_{inst} , as predicted by Theorem 3.

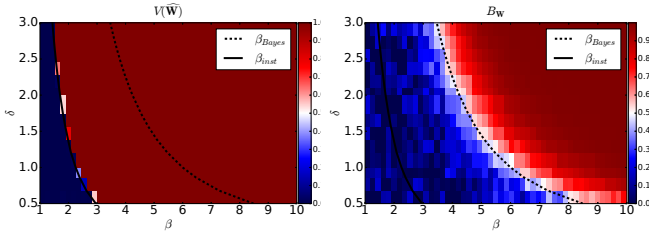


Figure 2. Gaussian model. Left frame: Empirical fraction of instances such that $V(\widehat{\mathbf{W}}) \geq \varepsilon_0 = 10^{-4}$, where $\widehat{\mathbf{W}}$ is the naive mean field estimate. Here $k = 2$, $d = 1000$ and, for each (δ, β) point on a grid, we used 400 random realizations to estimate the probability of $V(\widehat{\mathbf{W}}) \geq \varepsilon_0$. Right frame: Binder cumulant for the correlation between the naive mean field estimate $\widehat{\mathbf{W}}$, and the true weights \mathbf{W} , \mathbf{H} . Here $k = 2$, $d = 1000$ and $n = d\delta$, and we averaged over 400 realizations. Solid lines: The instability curve $\beta_{\text{inst}}(\delta)$. Dashed lines: The Bayes phase transition β_{Bayes} .

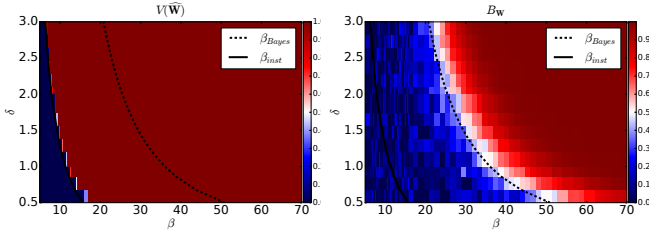


Figure 3. Same as for Figure 2, for the Dirichlet model.

Because of Proposition 3.1, we expect the estimates $\widehat{\mathbf{H}}$, $\widehat{\mathbf{W}}$ produced by variational inference to be asymptotically uncorrelated with the true factors for $\beta_{\text{inst}} < \beta < \beta_{\text{Bayes}}$. In order to test this hypothesis, we compute suitable correlation ratios between $\widehat{\mathbf{H}}$, $\widehat{\mathbf{W}}$ and the true parameters (known as ‘Binder cumulants’ $B_{\mathbf{H}}$ and $B_{\mathbf{W}}$). These quantities grow from 0 to 1 as β grows, and the transition is centered around β_{Bayes} . Figure 2 (right frame) reports the results $B_{\mathbf{W}}$ on a grid of (β, δ) values. Again, the transition is well predicted by the analytical curve β_{Bayes} . These data support our claim that, for $\beta_{\text{inst}} < \beta < \beta_{\text{Bayes}}$, the output of variational inference is non-uniform but uncorrelated with the true signal.

In Figure 3, we repeat the same experiment carried out in Figure 2, but for the Dirichlet model, with $\tilde{\nu} = \nu = 1$.

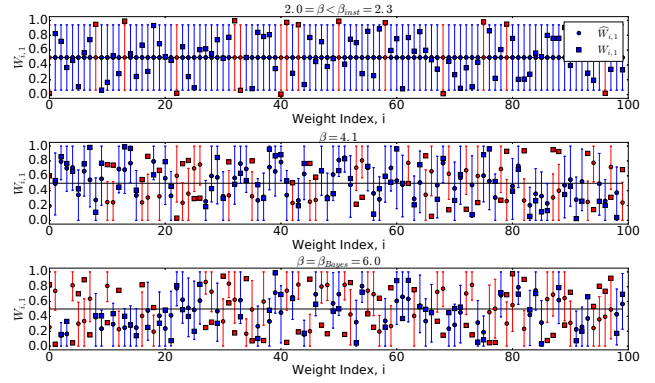


Figure 4. Gaussian model. Bayesian credible intervals as computed by variational inference at nominal coverage level $1 - \alpha = 0.9$. Here $k = 2$, $n = d = 5000$, $\beta \in \{2, 4.1, 6\}$ (for reference $\beta_{\text{inst}} \approx 2.3$, $\beta_{\text{Bayes}} = 6$). Circles: posterior mean. Squares: true weights. Red: coordinates on which the credible interval does not cover the true value of $w_{i,1}$.

Once more we observe a large region of model parameters for which the variational posterior is not centered on the uninformative point, but is uncorrelated with the ground truth. The instability phase transition is well captured by our theoretical prediction β_{inst} also in the Dirichlet case. Finally, in Figure 4 we plot the estimates obtained for 100 entries of the weights vector $w_{i,1}$ for three instances with $n = d = 5000$ and $\beta = 2 < \beta_{\text{inst}}$, $\beta = 4.1 \in (\beta_{\text{inst}}, \beta_{\text{Bayes}})$ and $\beta = 6 = \beta_{\text{Bayes}}$. The interval for $w_{a,1}$ is the form $\{w_{a,1} \in [0, 1] : \tilde{q}_a(w_{a,1}) \geq t_a(\alpha)\}$ and are constructed to achieve nominal coverage level $1 - \alpha = 0.9$. It is visually clear that the claimed coverage level is not verified in these simulations for $\beta > \beta_{\text{inst}}$, confirming our analytical results. Indeed, for the three simulations in Figure 4 we achieve coverage 0.87 (for $\beta = 2 < \beta_{\text{inst}}$), 0.65 (for $\beta = 4.1 \in (\beta_{\text{inst}}, \beta_{\text{Bayes}})$), and 0.51 (for $\beta = 6 = \beta_{\text{Bayes}}$). Further results of this type are reported in the SM.

4. Fixing the Instability

The fact that naive mean field is not accurate for certain classes of random high-dimensional probability distributions is well understood within statistical physics. In particular, in the context of mean field spin glasses (Mezard et al., 1988), naive mean field is known to lead to an asymptotically incorrect expression for the free energy. We expect the same mechanism to be relevant for topic models.

Namely, the product-form expression (1.3) only holds asymptotically in the sense of finite-dimensional marginals. However, when computing the term $\mathbb{E}_{\tilde{q}} \log p_{\mathbf{X}|\mathbf{W},\mathbf{H}}(\mathbf{X}|\mathbf{H},\mathbf{W})$ in the KL divergence (1.2), the error due to the product form approximation is non-negligible. Keeping track of this error leads to the so-called TAP free energy.

The TAP approach replaces the free energy (3.5) with

$\mathcal{F}_{\text{TAP}} = \mathcal{F}_{\text{TAP}}(\mathbf{r}, \tilde{\mathbf{r}})$ defined as follows (see SM)

$$\begin{aligned} \mathcal{F}_{\text{TAP}}(\mathbf{r}, \tilde{\mathbf{r}}) &= -\sqrt{\beta} \text{Tr}(\mathbf{X} \mathbf{r} \tilde{\mathbf{r}}^{\text{T}}) - \frac{\beta}{2d} \sum_{i=1}^d \sum_{a=1}^n \langle \mathbf{r}_i, \tilde{\mathbf{r}}_a \rangle^2 \\ &+ \sum_{i=1}^d \psi\left(\mathbf{r}_i, \frac{\beta}{d} \sum_{a=1}^n \tilde{\mathbf{r}}_a^{\otimes 2}\right) + \sum_{a=1}^n \tilde{\psi}\left(\tilde{\mathbf{r}}_a, \frac{\beta}{d} \sum_{i=1}^d \mathbf{r}_i^{\otimes 2}\right), \end{aligned}$$

where $\tilde{\mathbf{r}} \mathbf{1}_k = \mathbf{1}_n$, and we defined the partial Legendre transforms $\psi(\mathbf{r}, \mathbf{Q}) \equiv \sup_{\mathbf{m}} \{\langle \mathbf{r}, \mathbf{m} \rangle - \phi(\mathbf{m}, \mathbf{Q})\}$ and similarly for $\tilde{\psi}(\tilde{\mathbf{r}}, \tilde{\mathbf{Q}})$.

Calculus shows that stationary points of this free energy are in one-to-one correspondence with the fixed points of the following iteration:

$$\begin{aligned} \mathbf{m}^{t+1} &= \mathbf{X}^{\text{T}} \tilde{\mathbf{F}}(\tilde{\mathbf{m}}^t; \tilde{\mathbf{Q}}^t) - \mathbf{F}(\mathbf{m}^t; \mathbf{Q}^t) \tilde{\Omega}_t, \\ \tilde{\mathbf{m}}^t &= \mathbf{X} \mathbf{F}(\mathbf{m}^t; \mathbf{Q}^t) - \tilde{\mathbf{F}}(\tilde{\mathbf{m}}^{t-1}; \tilde{\mathbf{Q}}^{t-1}) \Omega_t, \\ \mathbf{Q}^{t+1} &= \frac{1}{d} \sum_{a=1}^n \tilde{\mathbf{F}}(\tilde{\mathbf{m}}_a^t; \tilde{\mathbf{Q}}^t)^{\otimes 2}, \quad \tilde{\mathbf{Q}}^t = \frac{1}{d} \sum_{i=1}^d \mathbf{F}(\mathbf{m}_i^t; \mathbf{Q}^t)^{\otimes 2}. \end{aligned}$$

where $\Omega_t, \tilde{\Omega}_t \in \mathbb{R}^{k \times k}$ are matrices defined in the SM. The stationarity conditions for the TAP free energy \mathcal{F}_{TAP} are known as TAP equations, and the above iterative algorithm is a special case of approximate message passing (AMP), with Bayesian updates.

Estimates of the factors \mathbf{W}, \mathbf{H} are computed following the same recipe as for naive mean field, cf. Eq. (??), namely $\hat{\mathbf{H}}^t = \mathbf{r}^t = \mathbf{F}(\mathbf{m}^t; \mathbf{Q}_t) / \sqrt{\beta}$, $\hat{\mathbf{W}}^t = \tilde{\mathbf{r}}^t = \tilde{\mathbf{F}}(\tilde{\mathbf{m}}^t; \tilde{\mathbf{Q}}_t) / \sqrt{\beta}$.

It is not hard to see that the AMP iteration admits an uninformative fixed point, which is a stationary point of the TAP free energy. This construction is analogous to the one for naive mean field, and we omit it here.

The next theorem establishes that the uninformative fixed point of the TAP free energy is a local minimum for all β below the spectral threshold $\beta_{\text{spect}}(k, \nu, \delta)$. Since $\beta_{\text{Bayes}}(k, \nu, \delta) \leq \beta_{\text{spect}}(k, \nu, \delta)$, this shows that the instability we discovered in the case of naive mean field is corrected by the TAP free energy.

Theorem 4. *Let $(\mathbf{r}_*, \tilde{\mathbf{r}}_*)$ be the uninformative stationary point of the TAP free energy. If $\beta < \beta_{\text{spect}}(k, \nu, \delta)$, then there exists $\varepsilon > 0$ such that, with high probability $\lambda_{\min}(\nabla^2 \mathcal{F}_{\text{TAP}}|_{(\mathbf{r}_*, \tilde{\mathbf{r}}_*)}) \geq \varepsilon$.*

In order to confirm the stability analysis at the previous section, we carried out numerical simulations analogous to the ones of Section 3.5. We initialize the iteration as for naive mean field, and monitor the same quantities, as in Section 3.5. In Figure 5 we report results on the distance from the uninformative subspace $V(\hat{\mathbf{W}})$, (left frame), and the Binder

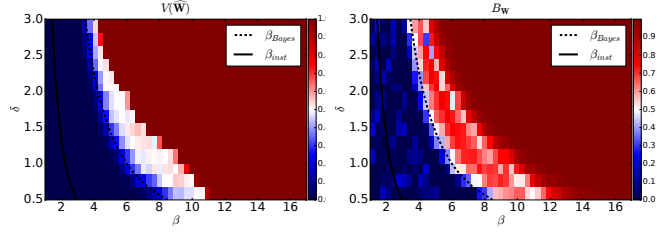


Figure 5. Left frame: Empirical fraction of instances such that $V(\hat{\mathbf{W}}) \geq \varepsilon_0 = 5 \cdot 10^{-3}$, where $\hat{\mathbf{W}}$ is the AMP estimate. Here $k = 2, d = 1000$, and for each (δ, β) point on the grid we ran AMP on 400 random realizations. Right frame: Binder cumulant for the correlation between AMP estimates $\hat{\mathbf{W}}, \hat{\mathbf{H}}$ and the true weights and topics \mathbf{W}, \mathbf{H} . Here $k = 2, d = 1000$ and estimates are obtained by averaging over 400 realizations.

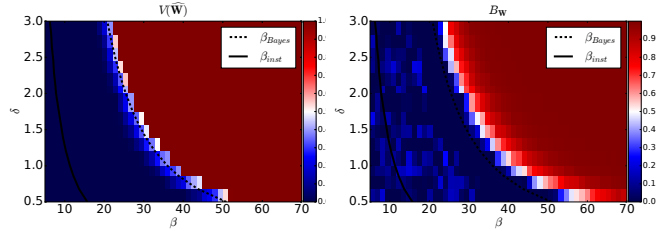


Figure 6. Same as in Figure 5, except for the Dirichlet model.

cumulant $B_{\mathbf{W}}$, measuring the correlation between AMP estimates and the true factors \mathbf{W}, \mathbf{H} (right frame). We repeat the same experiment in Figure 6 for the Dirichlet model.

In the intermediate regime $\beta \in (\beta_{\text{inst}}, \beta_{\text{spect}})$, the behavior of AMP is strikingly different from the one of naive mean field. AMP remains close to the uninformative fixed point, confirming that this is a local minimum of the TAP free energy. The distance from the uninformative subspace starts growing only at the spectral threshold β_{spect} (which coincides, in the present cases, with the Bayes threshold β_{Bayes}). At the same point, the correlation with the true factors \mathbf{W}, \mathbf{H} also becomes strictly positive.

5. Discussion

Bayesian methods are particularly attractive in unsupervised learning problems such as topic modeling. Even after a low-rank factorization $\mathbf{X} \approx \mathbf{W} \mathbf{H}^{\text{T}}$ is computed, it is still unclear how to evaluate it, or to which extent it should be trusted. Bayesian approaches provide estimates of the factors \mathbf{W}, \mathbf{H} , but also a probabilistic measure of how much these estimates should be trusted. To the extent that the posterior concentrates around its mean, this can be considered as a good estimate of a true underlying signal.

It is well understood that Bayesian estimates can be unreliable if the prior is not chosen carefully. Our work points at a second reason for caution. When variational inference is used for approximating the posterior, the result can be substantially incorrect.

Acknowledgements

H.J. and A.M. were partially supported by grants NSF CCF-1714305 and NSF IIS-1741162. B.G. was supported by Stanford’s Caroline and Fabian Pease Graduate Fellowship and grants NSF-DMS 1418362 and NSF-DMS 1407813.

References

- Airoldi, E. M., Blei, D. M., Fienberg, S. E., and Xing, E. P. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9(Sep):1981–2014, 2008.
- Alaoui, A. E. and Krzakala, F. Estimation in the spiked wigner model: A short proof of the replica formula. *arXiv preprint arXiv:1801.01593*, 2018.
- Arora, S., Ge, R., Kannan, R., and Moitra, A. Computing a nonnegative matrix factorization—provably. In *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*, pp. 145–162. ACM, 2012a.
- Arora, S., Ge, R., and Moitra, A. Learning topic models—going beyond svd. In *Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on*, pp. 1–10. IEEE, 2012b.
- Baik, J., Arous, G. B., Péché, S., et al. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *The Annals of Probability*, 33(5):1643–1697, 2005.
- Barbier, J., Dia, M., Macris, N., Krzakala, F., Lesieur, T., and Zdeborová, L. Mutual information for symmetric rank-one matrix estimation: A proof of the replica formula. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pp. 424–432. Curran Associates Inc., 2016.
- Bayati, M. and Montanari, A. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2):764–785, 2011.
- Benaych-Georges, F. and Nadakuditi, R. R. The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices. *Advances in Mathematics*, 227(1):494–521, 2011.
- Berthier, R., Montanari, A., and Nguyen, P.-M. State evolution for approximate message passing with non-separable functions. *arXiv preprint arXiv:1708.03950*, 2017.
- Bickel, P., Choi, D., Chang, X., Zhang, H., et al. Asymptotic normality of maximum likelihood and its variational approximation for stochastic blockmodels. *The Annals of Statistics*, 41(4):1922–1943, 2013.
- Binder, K. Finite size scaling analysis of ising model block distribution functions. *Zeitschrift für Physik B Condensed Matter*, 43(2):119–140, 1981.
- Blei, D. and Lafferty, J. Correlated topic models. *Advances in neural information processing systems*, 18:147, 2006a.
- Blei, D. M. Probabilistic topic models. *Communications of the ACM*, 55(4):77–84, 2012.
- Blei, D. M. and Lafferty, J. D. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning*, pp. 113–120. ACM, 2006b.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- Blei, D. M., Lafferty, J. D., et al. A correlated topic model of science. *The Annals of Applied Statistics*, 1(1):17–35, 2007.
- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- Broderick, T., Boyd, N., Wibisono, A., Wilson, A. C., and Jordan, M. I. Streaming variational bayes. In *Advances in Neural Information Processing Systems*, pp. 1727–1735, 2013.
- Celisse, A., Daudin, J.-J., Pierre, L., et al. Consistency of maximum-likelihood and variational estimators in the stochastic block model. *Electronic Journal of Statistics*, 6:1847–1899, 2012.
- Chang, J. and Blei, D. Relational topic models for document networks. In *Artificial Intelligence and Statistics*, pp. 81–88, 2009.
- Deshpande, Y. and Montanari, A. Information-theoretically optimal sparse pca. *arXiv preprint arXiv:1402.2238*, 2014.
- Deshpande, Y., Abbe, E., and Abbe, E. Asymptotic mutual information for the balanced binary stochastic block model. *Information and Inference: A Journal of the IMA*, 6(2):125–170, 2016.
- Donoho, D. L., Maleki, A., and Montanari, A. Message-passing algorithms for compressed sensing. *Proceedings of the National Academy of Sciences*, 106(45):18914–18919, 2009.
- Donoho, D. L., Maleki, A., and Montanari, A. Message passing algorithms for compressed sensing: I. motivation and construction. In *Information Theory (ITW 2010, Cairo), 2010 IEEE Information Theory Workshop on*, pp. 1–5. IEEE, 2010.

- Erosheva, E., Fienberg, S., and Lafferty, J. Mixed-membership models of scientific publications. *Proceedings of the National Academy of Sciences*, 101(suppl 1): 5220–5227, 2004.
- Fei-Fei, L. and Perona, P. A bayesian hierarchical model for learning natural scene categories. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 2, pp. 524–531. IEEE, 2005.
- Féral, D. and Pécché, S. The largest eigenvalue of rank one deformation of large wigner matrices. *Communications in mathematical physics*, 272(1):185–228, 2007.
- Giordano, R. J., Broderick, T., and Jordan, M. I. Linear response methods for accurate covariance estimates from mean field variational bayes. In *Advances in Neural Information Processing Systems*, pp. 1441–1449, 2015.
- Hoffman, M., Bach, F. R., and Blei, D. M. Online learning for latent dirichlet allocation. In *advances in neural information processing systems*, pp. 856–864, 2010.
- Javanmard, A. and Montanari, A. State evolution for general approximate message passing algorithms, with applications to spatial coupling. *Information and Inference: A Journal of the IMA*, 2(2):115–144, 2013.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- Korada, S. B. and Macris, N. Exact solution of the gauge symmetric p-spin glass model on a complete graph. *Journal of Statistical Physics*, 136(2):205–230, 2009.
- Krzakala, F., Xu, J., and Zdeborová, L. Mutual information in rank-one matrix estimation. *arXiv preprint arXiv:1603.08447*, 2016.
- Lelarge, M. and Miolane, L. Fundamental limits of symmetric low-rank matrix estimation. *Probability Theory and Related Fields*, pp. 1–71, 2016.
- Lesieur, T., Krzakala, F., and Zdeborová, L. Constrained low-rank matrix estimation: Phase transitions, approximate message passing and applications. *Journal of Statistical Mechanics: Theory and Experiment*, 2017(7): 073403, 2017.
- Mezard, M., Parisi, G., Virasoro, M. A., and Thouless, D. J. Spin glass theory and beyond. *Physics Today*, 41:109, 1988.
- Miolane, L. Fundamental limits of low-rank matrix estimation: the non-symmetric case. *arXiv preprint arXiv:1702.00473*, 2017.
- Montanari, A. and Venkataramanan, R. Estimation of low-rank matrices via approximate message passing. *arXiv preprint arXiv:1711.01682*, 2017.
- Opper, M. and Winther, O. Adaptive and self-averaging thouless-anderson-palmer mean-field theory for probabilistic modeling. *Physical Review E*, 64(5):056131, 2001.
- Pati, D., Bhattacharya, A., and Yang, Y. On statistical optimality of variational bayes. *arXiv preprint arXiv:1712.08983*, 2017.
- Perko, L. *Differential equations and dynamical systems*, volume 7. Springer Science & Business Media, 2013.
- Raj, A., Stephens, M., and Pritchard, J. K. Variational inference of population structure in large snp datasets. *Genetics*, pp. genetics–114, 2014.
- Recht, B., Re, C., Tropp, J., and Bittorf, V. Factoring non-negative matrices with linear programs. In *Advances in Neural Information Processing Systems*, pp. 1214–1222, 2012.
- Thouless, D. J., Anderson, P. W., and Palmer, R. G. Solution of ‘solvable model of a spin glass’. *Philosophical Magazine*, 35(3):593–601, 1977.
- Wainwright, M. J., Jordan, M. I., et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- Wang, B. and Titterton, D. Convergence and asymptotic normality of variational bayesian approximations for exponential family models with missing values. In *ARTIFICIAL INTELLIGENCE*. Citeseer, 2004.
- Wang, B. and Titterton, D. Inadequacy of interval estimates corresponding to variational bayesian approximations. In *AISTATS*. Barbados, 2005.
- Wang, B., Titterton, D., et al. Convergence properties of a general algorithm for calculating variational bayesian estimates for a normal mixture model. *Bayesian Analysis*, 1(3):625–650, 2006.
- Wang, C. and Blei, D. M. Collaborative topic modeling for recommending scientific articles. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 448–456. ACM, 2011.
- Wang, X. and Grimson, E. Spatial latent dirichlet allocation. In *Advances in neural information processing systems*, pp. 1577–1584, 2008.

Zhang, A. Y. and Zhou, H. H. Theoretical and computational guarantees of mean field variational inference for community detection. *arXiv preprint arXiv:1710.11268*, 2017.

Zhou, J., Bhattacharya, A., Herring, A. H., and Dunson, D. B. Bayesian factorizations of big sparse tensors. *Journal of the American Statistical Association*, 110(512): 1562–1576, 2015.