

Consistency of ELBO maximization for model selection

Badr-Eddine Chérif-Abdellatif
CREST, ENSAE, Université Paris Saclay

BADR.EDDINE.CHERIEF.ABDELLATIF@ENSAE.FR

Abstract

The Evidence Lower Bound (ELBO) is a quantity that plays a key role in variational inference. It can also be used as a criterion in model selection. However, though extremely popular in practice in the variational Bayes community, there has never been a general theoretic justification for selecting based on the ELBO. In this paper, we show that the ELBO maximization strategy has strong theoretical guarantees, and is robust to model misspecification while most works rely on the assumption that one model is correctly specified. We illustrate our theoretical results by an application to the selection of the number of principal components in probabilistic PCA.

Keywords: Variational inference, Evidence lower bound, Model selection.

1. Introduction

Approximate Bayesian inference is at the core of modern Bayesian statistics and machine learning. While exact Bayesian inference is often intractable, variational inference has proved to provide an efficient solution when dealing with large datasets and complex probabilistic models. Variational Bayes (VB) aims at maximizing a numerical quantity referred to as Evidence Lower Bound on the marginal likelihood (ELBO), and thus makes use of optimization techniques to converge faster than Monte Carlo sampling approach. [Blei et al. \(2017\)](#) provides a comprehensive survey on variational inference. Although VB is mainly used for its practical efficiency, little attention has been put towards its theoretical properties during the last years. While [Alquier et al. \(2016\)](#) studied the properties of variational approximations of Gibbs distributions used in machine learning for bounded loss functions, [Alquier and Ridgway \(2017\)](#); [Zhang and Gao \(2017\)](#); [Wang and Blei \(2018\)](#); [Bhattacharya et al. \(2018\)](#); [Chérif-Abdellatif and Alquier \(2018\)](#) extended the results to more general statistical models.

At the same time, model selection remains a major problem of interest in statistics that naturally arises in the course of scientific inquiry. The statistician aims at selecting a model among several candidates given an observed dataset. To do so, one can perform cross validation as in [Vehtari et al. \(2014\)](#) or maximize a numerical criterion to make the final choice, see the review of [Rao and Wu \(2001\)](#). In the literature, penalized criteria such as AIC and BIC respectively introduced by [Akaike \(1974\)](#) and [Schwarz \(1978\)](#) are popular. While AIC aims at optimizing the prediction performance, BIC is more suitable for recovering with high probability the true model (when such a model exists), see [Yang \(2005\)](#). Thus, it is necessary to define a criterion suited to a given objective. Meanwhile, a non-asymptotic theory of penalization using oracle inequalities has been developed during

the last two decades, and offers a simple way to assess the quality of a given model selection criterion. We refer the interested reader to [Massart \(2007\)](#) for more details.

In this paper, we are interested in finding an estimate of the distribution of the data, and we need to choose from among competing models. [Blei et al. \(2017\)](#) states that "the [evidence lower] bound is a good approximation of the marginal likelihood, which provides a basis for selecting a model. Though this sometimes works in practice, selecting based on a bound is not justified in theory". Since then, authors of [Chérif-Abdellatif and Alquier \(2018\)](#) have provided an analysis of model selection based on the ELBO in the case of mixture models. We extend their result to the general case of independent and identically distributed (i.i.d.) data, and we provide an oracle inequality on the ELBO criterion that justifies the consistency of ELBO maximization when the objective is the estimation of the distribution of the data. In particular, as soon as there exists a true model, we show that the ELBO criterion is adaptive and that the selected estimator achieves the convergence rate of the variational approximation associated with the true model.

The rest of this paper is organized as follows. Section 2 introduces the setting and the key concepts needed to understand our results. In Section 3, we prove that the ELBO criterion provides a variational approximation that is consistent with the sample size as soon as there exists a true model. We also extend the result to misspecified models. We finally illustrate the main theorem of this paper by an application to the selection of the number of principal components in probabilistic Principal Component Analysis (PCA) in Section 4. All the proofs are deferred to the appendix.

2. Framework

Let us introduce the notations and the framework we adopt in this paper. We consider a collection of i.i.d. random variables X_1, \dots, X_n distributed according to some probability distribution P^0 in a measurable space $(\mathbb{X}, \mathcal{X})$. We denote $X_1^n = (X_1, \dots, X_n)$. We consider a countable collection $\{\mathcal{M}_K / K \geq 1\}$ of statistical mixture models $\mathcal{M}_K = \{P_{\theta_K} / \theta_K \in \Theta_K\}$ where Θ_K is the parameter set associated with index K . We make no assumptions on Θ_K 's nor on P_{θ_K} . Parameter spaces may overlap or have inclusion relationships. Let $\mathcal{M}_1^+(\Theta_K)$ be the set of all probability distributions over Θ_K .

We use a Bayesian approach, and we define a prior π over the full parameter space $\cup_{K \geq 1} \Theta_K$ (equipped with some suited sigma-algebra). First, we specify a prior weight π_K assigned to model \mathcal{M}_K , and then a conditional prior $\Pi_K(\cdot)$ on $\theta_K \in \Theta_K$ given model \mathcal{M}_K :

$$\pi = \sum_{K \geq 1} \pi_K \Pi_K.$$

The Kullback-Leibler divergence between two probability distributions P and R is

$$\text{KL}(P, R) = \begin{cases} \int \log \left(\frac{dP}{dR} \right) dP & \text{if } R \text{ dominates } P, \\ +\infty & \text{otherwise.} \end{cases}$$

For any $\alpha \neq 1$, authors of [Van Erven and Harremos \(2014\)](#) detail the properties of the α -Rényi divergence between two probability distributions P and R which is equal to:

$$D_\alpha(P, R) = \begin{cases} \frac{1}{\alpha-1} \log \int \left(\frac{dP}{dR} \right)^{\alpha-1} dP & \text{if } R \text{ dominates } P, \\ +\infty & \text{otherwise.} \end{cases}$$

We define the tempered posterior distribution $\pi_{n,\alpha}^K(\cdot|X_1^n)$ on parameter $\theta_K \in \Theta_K$ given model \mathcal{M}_K using prior Π_K and likelihood L_n for any $\alpha \in (0, 1)$:

$$\pi_{n,\alpha}^K(d\theta_K|X_1^n) \propto L_n(\theta_K)^\alpha \Pi_K(d\theta_K).$$

This definition is a slight variant of the regular Bayesian posterior (for which $\alpha = 1$), and is also referred to as Bayesian fractional posterior in [Bhattacharya et al. \(2016\)](#). This posterior is easier to sample from, more robust to model misspecification and requires less stringent conditions to obtain consistency, see respectively [Behrens et al. \(2012\)](#), [Grünwald and Van Ommen \(2017\)](#) and [Bhattacharya et al. \(2016\)](#).

The Variational Bayes approximation $\tilde{\pi}_{n,\alpha}^K(\cdot|X_1^n)$ of the tempered posterior associated with model \mathcal{M}_K is then defined as the projection, with respect to the Kullback-Leibler divergence, of the tempered posterior onto some set \mathcal{F}_K :

$$\tilde{\pi}_{n,\alpha}^K(\cdot|X_1^n) = \arg \min_{\rho_K \in \mathcal{F}_K} \text{KL}(\rho_K, \pi_{n,\alpha}^K(\cdot|X_1^n)).$$

The choice of the variational set \mathcal{F}_K is crucial: the variational approximation must be close enough to the target distribution (as an approximation of the tempered posterior) but not too close (in order to be tractable). A classical variational set \mathcal{F}_K is the parametric family which leads to a tractable parametric approximation, e.g. a Gaussian distribution. Another popular set \mathcal{F}_K in the VB community is the mean-field approximation that is based on a partition of the space of parameters, and which consists in a factorization of the variational approximation over the partition.

Alternatively, the variational approximation is often defined as the distribution into \mathcal{F}_K that maximizes the Evidence Lower Bound:

$$\tilde{\pi}_{n,\alpha}^K(\cdot|X_1^n) = \arg \max_{\rho_K \in \mathcal{F}_K} \left\{ \alpha \int \ell_n(\theta_K) \rho_K(d\theta_K) - \text{KL}(\rho_K, \Pi_K) \right\}$$

where the function inside the argmax operator is the ELBO (as a function of K and ρ_K) and ℓ_n is the log-likelihood. In the following, we will just call $\text{ELBO}(K)$ the closest approximation to the log-evidence, i.e. the value of the ELBO evaluated at its maximum:

$$\text{ELBO}(K) = \alpha \int \ell_n(\theta_K) \tilde{\pi}_{n,\alpha}^K(d\theta_K|X_1^n) - \text{KL}(\tilde{\pi}_{n,\alpha}^K(\cdot|X_1^n), \Pi_K).$$

In the variational Bayes community, researchers and practitioners use the ELBO in order to select the model from which they will consider the final variational approximation $\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n)$, as stated in [Blei et al. \(2017\)](#). We propose to consider a penalized version of the ELBO criterion

$$\hat{K} = \arg \max_{K \geq 1} \left\{ \text{ELBO}(K) - \log \left(\frac{1}{\pi_K} \right) \right\}$$

which is a slight variant of the classical definition, although choosing a uniform prior over a finite number of models leads to maximizing the ELBO. Note that the penalty term is not just an artefact in order to ease the theoretical proof, but it is a complexity term that reflects our prior beliefs over the different models.

We will provide in the next section a theoretical justification to such a selection criterion and show that the selected variational estimator $\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n)$ is consistent under mild conditions as soon as there exists a true model. We will adopt the definition of *consistency* used in [Alquier and Ridgway \(2017\)](#) and [Chérif-Abdellatif and Alquier \(2018\)](#) that is, the Bayesian estimator is said to be consistent if, in expectation (with respect to the random variables distributed according to P^0), the average Rényi loss between a distribution in the selected model and the true distribution (over the Bayesian estimator) goes to zero as $n \rightarrow +\infty$:

$$\mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta|X_1^n) \right] \xrightarrow{n \rightarrow +\infty} 0.$$

This definition is closely related to the notion of *concentration* which is defined in [Ghosal et al. \(2000\)](#) as the asymptotic concentration of the Bayesian estimator around the true distribution, and which is usually used to assess frequentist guarantees for Bayesian estimators. It is sometimes also referred to as *contraction* (or even *consistency*). See [Appendix A](#) for more details on the connection between the notions of *consistency* and *concentration*.

3. Consistency of the ELBO criterion

In this section, unless explicitly stated otherwise, we assume that there exists a true model \mathcal{M}_{K_0} that contains the true distribution P^0 , i.e. that there exists K_0 and $\theta^0 \in \Theta_{K_0}$ such that $P^0 = P_{\theta^0}$.

A key assumption introduced in [Ghosal et al. \(2000\)](#) in order to obtain the concentration of the regular posterior distribution $\pi_{n,1}^{K_0}(\cdot|X_1^n)$ associated with the true model \mathcal{M}_{K_0} is a *prior mass condition* which states that the prior Π_{K_0} must give enough mass to some neighborhood (in the Kullback-Leibler sense) of the true parameter. [Bhattacharya et al. \(2016\)](#) showed that this condition was sufficient when considering tempered posteriors $\pi_{n,\alpha}^{K_0}(\cdot|X_1^n)$. [Alquier and Ridgway \(2017\)](#) extended this assumption in order to obtain the concentration and the consistency of variational approximations of the tempered posteriors $\tilde{\pi}_{n,\alpha}^{K_0}(\cdot|X_1^n)$. In addition to the previous prior mass condition, this extension requires the variational set \mathcal{F}_{K_0} to contain probability distributions concentrated around the true parameter. Note that when $\mathcal{F}_{K_0} = \mathcal{M}_1^+(\Theta_{K_0})$, this goes back to the standard prior mass condition. This extended prior mass condition is standard in the variational Bayes community, see [Alquier and Ridgway \(2017\)](#); [Chérif-Abdellatif and Alquier \(2018\)](#), and can be formulated as follows:

Assumption : *We assume that there exists r_n for which there is a distribution $\rho_{K_0,n} \in \mathcal{F}_{K_0}$ such that:*

$$\int \text{KL}(P^0, P_{\theta_{K_0}}) \rho_{K_0,n}(d\theta_{K_0}) \leq r_n \text{ and } \text{KL}(\rho_{K_0,n}, \Pi_{K_0}) \leq nr_n. \quad (3.1)$$

Remark 1 *Define the KL-ball \mathcal{B} centered at θ_0 of radius r_n :*

$$\mathcal{B} = \{\theta \in \Theta_{K_0} / \text{KL}(P_{\theta_0}, P_\theta) \leq r_n\},$$

and consider the restriction $\rho_{K_0,n}$ of Π_{K_0} to \mathcal{B} . Then it is clear that when $\rho_{K_0,n} \in \mathcal{F}_{K_0}$, [Assumption 3.1](#) becomes equivalent to the former prior mass condition of [Ghosal et al. \(2000\)](#), i.e. $\Pi_{K_0}(\mathcal{B}) \geq e^{-nr_n}$. The computation of the prior mass $\Pi_{K_0}(\mathcal{B})$ is a major difficulty.

It has been raised as a question of interest in Ghosal et al. (2000), and is addressed for categorical distributions and Dirichlet priors in Ghosal et al. (2000) (but for an L_1 -ball) and in Chérief-Abdellatif and Alquier (2018) (for a KL-ball). Unfortunately, $\rho_{K_0,n}$ does not belong to \mathcal{F}_{K_0} in general and the computation of the prior mass is no longer sufficient. Nevertheless, the strategy of computing the prior mass of KL-balls remains of interest when dealing with mixture models and mean-field approximation sets, see Chérief-Abdellatif and Alquier (2018) where the authors showed that studying the prior mass condition of Ghosal et al. (2000) independently on the weights and on each component becomes sufficient.

Remark 2 When \mathcal{F}_{K_0} is parametric, it is often possible to overcome the difficulty presented above in order to find a rate r_n as in Assumption 3.1. Indeed, the point is to express the distribution $\rho_{K_0,n}$ using the general parametric form of the variational family, and to find relevant values of the parameters that will lead to fast rates of convergence r_n . This is the strategy we follow in Section 4 for probabilistic PCA. See Alquier and Ridgway (2017); Chérief-Abdellatif and Alquier (2018) for other examples of such computations.

Alquier and Ridgway (2017) showed that the variational approximation $\tilde{\pi}_{n,\alpha}^{K_0}(\cdot|X_1^n)$ associated with a true model is consistent under Assumption 3.1 and that the convergence rate is equal to r_n . Nevertheless, in model selection, we do not necessarily know which model is true and the challenge is to be able to find one such that the corresponding approximation is consistent at a comparable convergence rate. We show that the variational approximation $\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n)$ associated with the selected model is also consistent at rate r_n as soon as Assumption 3.1 is satisfied:

Theorem 3 Assume that Assumption 3.1 is satisfied. Then for any $\alpha \in (0, 1)$,

$$\mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta|X_1^n) \right] \leq \frac{1+\alpha}{1-\alpha} r_n + \frac{\log(\frac{1}{\pi_{K_0}})}{n(1-\alpha)}.$$

The inequality in Theorem 3 shows the adaptivity of our procedure. Indeed, whatever the value of \hat{K} (which can be different from K_0), we obtain the consistency of the selected variational approximation at the same rate of convergence than the estimator associated with the true model (as soon as the additional term in the upper bound is lower than r_n , which is the case for prior weights used in practice). We recall that we look for a good estimation of the true distribution P^0 and not for an estimation of the true model index K_0 which is a different task that would require identifiability assumptions that are stronger than those in our theorem. The overall rate is composed of the convergence rate associated with the true model \mathcal{M}_{K_0} , and of a complexity term that reflects the prior belief over the (unknown) true model. For example, if we range a countable number of models according to our prior belief, and we take $\pi_K = 2^{-K}$, then the corresponding term will be of order K_0/n . More generally, when $\frac{1}{n} \lesssim r_n$, we obtain the consistency at the rate associated with the true model.

As a short example, Chérief-Abdellatif and Alquier (2018) investigated the case of mixture models. For instance, authors obtained a convergence rate equal to $K_0 \log(nK_0)/n$ for Gaussian mixtures when there exists a true K_0 -components mixture model. We study another example in Section 4.

We can also extend this result to misspecified models. In the model selection literature, only little attention has been put to misspecification when the true distribution does not belong to any of the models, see [Lv and Liu \(2013\)](#). Now, we do not assume any longer that there exists a true model, and we show that our ELBO criterion is robust to model misspecification:

Theorem 4 *For each index K , let us define the set $\Theta_K(r_{K,n})$ of parameters $\theta_K^* \in \Theta_K$, for which there is a distribution $\rho_{K,n} \in \mathcal{F}_K$ such that:*

$$\int \mathbb{E} \left[\log \frac{P_{\theta_K^*}(X_i)}{P_{\theta_K}(X_i)} \right] \rho_{K,n}(d\theta_K) \leq r_{K,n} \text{ and } \text{KL}(\rho_{K,n}, \Pi_K) \leq nr_{K,n}. \quad (3.2)$$

Then for any $\alpha \in (0, 1)$,

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \\ & \leq \inf_{K \geq 1} \left\{ \frac{\alpha}{1 - \alpha} \inf_{\theta_K^* \in \Theta_K(r_{K,n})} \text{KL}(P^0, P_{\theta_K^*}) + \frac{1 + \alpha}{1 - \alpha} r_{K,n} + \frac{\log(\frac{1}{\pi_K})}{n(1 - \alpha)} \right\}. \end{aligned}$$

Note that when there exists a true model \mathcal{M}_{K_0} such that $P^0 = P_{\theta^0}$ with $\theta^0 \in \Theta_{K_0}$, then under Assumption 3.1, we get $\theta^0 \in \Theta_{K_0}(r_{K_0,n})$, and we recover Theorem 3. Furthermore, the oracle inequality in Theorem 4 shows that the selected variational approximation adaptively achieves the best upper bound among the different models \mathcal{M}_K , where each upper bound is a trade-off between two terms: a bias due to the error of approximating the true distribution by a distribution in model \mathcal{M}_K , and a variance term $r_{K,n}$ (as soon as the penalty term is lower than $r_{K,n}$) that is defined in Condition 3.2.

4. Application to probabilistic PCA

We consider here the probabilistic Principal Component Analysis (PCA) problem as an application of our work. From now on, matrices will be denoted in bold capital letters. We assume the model

$$X_i = \mathbf{W}Z_i + \sigma^2 \mathbf{I}_d$$

with i.i.d. Gaussian random variables $Z_i \sim \mathcal{N}(0, \mathbf{I}_K)$, where \mathbf{I}_d and \mathbf{I}_K are respectively the d - and K -dimensional identity matrices ($K < d$), $\mathbf{W} \in \mathbb{R}^{d \times K}$ is the K -rank matrix that contains the principal axes and σ^2 is a noisy term that is known. We suppose here that data are centred. Hence, the distribution of each X_i is

$$P_{\mathbf{W}} := \mathcal{N}(0, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d).$$

We are not interested here in estimating the principal axes \mathbf{W} and selecting the number of components K , but in estimating the true distribution of the X_i 's.

Each model corresponds to a rank K . We place an equal prior weight over each integer $K = 1, \dots, d$. Hence the optimization problem is equivalent to maximizing the ELBO as in [Blei et al. \(2017\)](#). Given rank K , we place a prior over the K -rank matrix \mathbf{W} to infer a distribution over principal axes. We choose independent Gaussian priors $\mathcal{N}(0, s^2 \mathbf{I}_d)$

on the columns W_1, \dots, W_K of \mathbf{W} . We also consider Gaussian independent variational approximations $\mathcal{N}(\mu_j, \Sigma_j)$ for the columns of \mathbf{W} . Then, as soon as there exists a true model, i.e. there exists K_0 and $\mathbf{W}_0 \in \mathbb{R}^{d \times K_0}$ such that the true distribution of each X_i is $P_{\mathbf{W}_0} = \mathcal{N}(0, \mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d)$, under the assumption that the coefficients of \mathbf{W}_0 are bounded, then Theorem 5 provides an explicit rate of convergence of our variational estimator even when K_0 is unknown:

Theorem 5 *For any $\alpha \in (0, 1)$, as soon as there exists a true model \mathcal{M}_{K_0} such that $P^0 = P_{\mathbf{W}_0}$ with $\mathbf{W}_0 \in \mathbb{R}^{d \times K_0}$ and such that the coefficients of \mathbf{W}_0 are bounded, then:*

$$\mathbb{E} \left[\int D_\alpha(P_{\mathbf{W}}, P_{\mathbf{W}_0}) \tilde{\pi}_{n, \alpha}^{\hat{K}}(d\mathbf{W} | X_1^n) \right] = \mathcal{O} \left(\frac{dK_0 \log(dn)}{n} \right).$$

The proof as well as the computation of the ELBO are detailed in the appendix. Note that this corollary can directly lead to a result in Frobenius distance between covariance matrices $\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d$ and $\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d$ instead of the Rényi divergence between the corresponding distributions even when \mathbf{W} and \mathbf{W}_0 are not equal-sized matrices. We denote $\|\cdot\|_F$ the Frobenius norm and $\|\cdot\|_2$ the spectral norm of a matrix, which are respectively defined as the square root of the sum of the absolute squares of the elements of a matrix and as its largest singular value.

The following corollary assesses the consistency of the selected variational approximation to the true covariance matrix in Frobenius norm. The idea, borrowed from [Alquier and Ridgway \(2017\)](#), is to project matrices onto some set of bounded matrices under the assumption that the spectral norm of the true matrix \mathbf{W}_0 is also bounded:

Corollary 6 *For any $\alpha \in (0, 1)$, as soon as there exists a true model \mathcal{M}_{K_0} such that $P^0 = P_{\mathbf{W}_0}$ with $\mathbf{W}_0 \in \mathbb{R}^{d \times K_0}$ and such that the spectral norm of \mathbf{W}_0 is upper bounded by a positive constant $B > 0$, then:*

$$\mathbb{E} \left[\int \|\text{clip}_B(\mathbf{W}\mathbf{W}^T) - \mathbf{W}_0 \mathbf{W}_0^T\|_F^2 \tilde{\pi}_{n, \alpha}^{\hat{K}}(d\mathbf{W} | X_1^n) \right] = \mathcal{O} \left(\frac{dK_0 \log(dn)}{n} \right)$$

where $\text{clip}_B(\mathbf{A})$ is the matrix which (i, j) -entry is equal to $\begin{cases} \mathbf{A}_{i,j} & \text{if } |\mathbf{A}_{i,j}| \leq B^2 \\ B^2 & \text{if } \mathbf{A}_{i,j} \geq B^2 \\ -B^2 & \text{otherwise.} \end{cases}$

The requirement in our corollary is that the spectral norm of the true matrix \mathbf{W}_0 is bounded by some positive constant B , which implies the boundedness of the coefficients of the matrix as required in Theorem 5. In particular, the coefficients of the matrix $\mathbf{W}_0 \mathbf{W}_0^T$ are bounded by B^2 :

$$\begin{aligned} |(\mathbf{W}_0 \mathbf{W}_0^T)_{i,j}| &= \left| \sum_{k=1}^{K_0} (\mathbf{W}_0)_{i,k} (\mathbf{W}_0)_{j,k} \right| \leq \left(\sum_{k=1}^{K_0} (\mathbf{W}_0)_{i,k}^2 \right)^{1/2} \left(\sum_{k=1}^{K_0} (\mathbf{W}_0)_{j,k}^2 \right)^{1/2} \\ &= \frac{\|\mathbf{W}_0 e_i\|_2}{\|e_i\|_2} \frac{\|\mathbf{W}_0 e_j\|_2}{\|e_j\|_2} \leq \|\mathbf{W}_0\|_2^2 \leq B^2 \end{aligned}$$

using Cauchy-Schwarz inequality and the property $\|\mathbf{W}_0\|_2 = \max_{x \neq 0} \frac{\|\mathbf{W}_0 x\|_2}{\|x\|_2}$ where e_ℓ is the vector of \mathbb{R}^d which components are all equal to 0 except for the ℓ -th one that is set to 1. Hence it seems sensible to project (with respect to the Frobenius distance) any estimator $\mathbf{W}\mathbf{W}^T$ onto the set of all matrices whose entries lie in the interval $[-B^2, B^2]$, which is exactly what the $clip_B$ application does. Note that the spectral norm of Matrix \mathbf{W}_0 is equal to the largest eigenvalue of $\mathbf{W}_0 \mathbf{W}_0^T$, so our assumption comes back to upper bounding the eigenvalues of the covariance matrix $\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d$, which is a classical assumption when estimating covariance matrices, see for instance [Cai et al. \(2015\)](#).

It is also possible to obtain a consistent pointwise covariance matrix estimator with the same convergence rate:

Corollary 7 *For any $\alpha \in (0, 1)$, as soon as there exists a true model \mathcal{M}_{K_0} such that $P^0 = P_{\mathbf{W}_0}$ with $\mathbf{W}_0 \in \mathbb{R}^{d \times K_0}$ and such that the spectral norm of \mathbf{W}_0 is bounded by B . Let us define a pointwise estimator of the covariance matrix:*

$$\hat{\Sigma} = \int clip_B(\mathbf{W}\mathbf{W}^T) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\mathbf{W}|X_1^n) + \sigma^2 \mathbf{I}_d.$$

Then,

$$\mathbb{E} \left[\left\| \hat{\Sigma} - (\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d) \right\|_F^2 \right] = \mathcal{O} \left(\frac{dK_0 \log(dn)}{n} \right).$$

Discussion

In this paper we proved the consistency of ELBO maximization in model selection. By penalizing the variational lower bound using our prior beliefs over the different models, we showed that under mild conditions, the variational approximation associated with the selected model is consistent at the same convergence rate than the approximation associated with the true model. Moreover, the oracle inequality in [Theorem 4](#) proved that the selected approximation is robust to misspecification. An application to the selection of the number of principal components in probabilistic PCA was provided as a short example.

We discuss in [Appendix A](#) the connection between the notions of *consistency* and *concentration*. This justifies the use of the α parameter in the definition of the evidence lower bound, as the regular posterior distribution is not robust to model misspecification. Indeed, authors of [Grünwald and Van Ommen \(2017\)](#) explain that there are pathologic cases where the regular posterior does not concentrate to the true distribution.

A point of interest when dealing with model selection is the question of recovering the true model (when it exists). This issue falls beyond the scope of this paper which treats the question of estimating the true distribution, and can be the object of future works. The true model recovery would require stronger assumptions, but the implementation in [Section 5](#) in [Bishop \(1999\)](#) suggests that those may hold for probabilistic PCA.

Also, it would be interesting to study cross-validation instead of ELBO maximization. However, the tools used in this work such as the theory of penalized criteria and oracle inequalities were particularly suited to the ELBO, and thus a different theory should be used in order to obtain the consistency of validation log-likelihood in the VB framework. This question is left for future research.

Acknowledgments

I would like to warmly thank Pierre Alquier, Lionel Riou-Durand and the anonymous referees for their inspiring comments and suggestions on this work.

References

- H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19:716–723, 1974.
- P. Alquier and J. Ridgway. Concentration of tempered posteriors and of their variational approximations. *arXiv preprint arXiv:1706.09293*, 2017.
- P. Alquier, J. Ridgway, and N. Chopin. On the properties of variational approximations of Gibbs posteriors. *JMLR*, 17(239):1–41, 2016.
- G. Behrens, N. Friel, and M. Hurn. Tuning tempered transitions. *Statistics and computing*, 22(1):65–78, 2012.
- A. Bhattacharya, D. Pati, and Y. Yang. Bayesian fractional posteriors. *arXiv preprint arXiv:1611.01125, to appear in the Annals of Statistics*, 2016.
- A. Bhattacharya, D. Pati, and Y. Yang. On statistical optimality of variational Bayes. *Proceedings of Machine Learning Research*, 84 - AISTAT, 2018.
- C. Bishop. Variational principal components. In *Proceedings Ninth International Conference on Artificial Neural Networks, ICANN’99*, volume 1, pages 509–514. IEE, January 1999.
- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- T. Cai, Z. Ma, and Y. Wu. Optimal estimation and rank detection for sparse spiked covariance matrices. *Probability Theory and Related Fields*, 161(3):781–815, Apr 2015. ISSN 1432-2064. doi: 10.1007/s00440-014-0562-z.
- O. Catoni. *PAC-Bayesian supervised classification: the thermodynamics of statistical learning*. Institute of Mathematical Statistics Lecture Notes—Monograph Series, 56. Institute of Mathematical Statistics, Beachwood, OH, 2007.
- B. Chérief-Abdellatif and P. Alquier. Consistency of variational bayes inference for estimation and model selection in mixtures. *Electronic Journal of Statistics*, 12(2):2995–3035, 2018. ISSN 1935-7524. doi: 10.1214/18-EJS1475.
- S. Forth, P. Hovland, E. Phipps, J. Utke, and A. Walther. *Recent Advances in Algorithmic Differentiation*. Springer Publishing Company, Incorporated, 2014. ISBN 3642439918, 9783642439919.
- S. Ghosal, J. K. Ghosh, and A. W. Van Der Vaart. Convergence rates of posterior distributions. *Annals of Statistics*, pages 500–531, 2000.

- P. D. Grünwald and T. Van Ommen. Inconsistency of Bayesian inference for misspecified linear models, and a proposal for repairing it. *Bayesian Analysis*, 12(4):1069–1103, 2017.
- J. Lv and J. S. Liu. Model selection principles in misspecified models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):141–167, 2013. doi: 10.1111/rssb.12023.
- P. Massart. *Concentration inequalities and model selection*, volume 1896 of *Lecture Notes in Mathematics*. Springer, Berlin, 2007. ISBN 978-3-540-48497-4; 3-540-48497-3. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003, Edited by Jean Picard.
- K.I. Moridomi, K. Hatano, and E. Takimoto. Online linear optimization with the log-determinant regularizer. *IEICE Transactions on Information and Systems*, E101D(6): 1511–1520, 6 2018. ISSN 0916-8532. doi: 10.1587/transinf.2017EDP7317.
- C. R. Rao and Y. Wu. *On model selection*, volume Volume 38 of *Lecture Notes–Monograph Series*, pages 1–57. Institute of Mathematical Statistics, Beachwood, OH, 2001. doi: 10.1214/lnms/1215540960.
- G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- T. Van Erven and P. Harremoës. Rényi divergence and Kullback-Leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820, 2014.
- A. Vehtari, V. Tolvanen, T. Mononen, and O. Winther. Bayesian leave-one-out cross validation approximations for gaussian latent variable models. *Journal of Machine Learning Research*, 17, 12 2014.
- Y. Wang and D. M. Blei. Frequentist consistency of variational Bayes. *Journal of the American Statistical Association*, pages 1–85, 2018.
- Y. Yang. Can the strengths of AIC and BIC be shared? A conflict between model identification and regression estimation. *Biometrika*, 92(4):937–950, 2005.
- F. Zhang and C. Gao. Convergence rates of variational posterior distributions. *arXiv preprint arXiv:1712.02519v1*, 2017.

Appendix A. Connection between consistency and concentration.

In this appendix, we highlight the connection between the notions of *consistency* used in [Alquier and Ridgway \(2017\)](#) and [Chérif-Abdellatif and Alquier \(2018\)](#) and *concentration*. We consider a true model \mathcal{M}_{K_0} to which the true distribution $P^0 = P_{\theta^0}$ belongs, $\theta^0 \in \Theta_{K_0}$. We recall that the Bayesian estimator $\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n)$ is said to be consistent if, in expectation (with respect to the random variables distributed according to P^0), the average Rényi loss between a distribution in the selected model and the true distribution (over the Bayesian estimator) goes to zero as $n \rightarrow +\infty$:

$$\mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta|X_1^n) \right] \xrightarrow{n \rightarrow +\infty} 0.$$

Similarly, we give the definition of *concentration* at rate s_n of the selected variational approximation to P^0 as stated in [Ghosal et al. \(2000\)](#), that is, in probability (with respect to the random variables distributed according to P^0), the approximation concentrates asymptotically around the true distribution as $n \rightarrow +\infty$, i.e. in probability:

$$\tilde{\pi}_{n,\alpha}^{\hat{K}} \left(D_\alpha(P_\theta, P^0) > Ms_n | X_1^n \right) \xrightarrow{n \rightarrow +\infty} 0$$

for any constant $M > 0$. The reference metric here is the α -Rényi divergence.

We show in this appendix that the consistency of the selected variational approximation to P^0 at rate r_n implies the concentration of the selected variational approximation to P^0 at any rate s_n such that $r_n = o(s_n)$ and $s_n \rightarrow 0$ as $n \rightarrow +\infty$, as for instance $s_n = r_n \log(\log(n))$ when the consistency rate r_n is slower than a log-logarithmic one.

To do so, we assume that the selected variational approximation is consistent to P^0 at rate r_n , i.e.:

$$\mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta|X_1^n) \right] \leq r_n.$$

Then, using Markov's inequality for any s_n such that $r_n = o(s_n)$ and $s_n \rightarrow 0$ and any constant $M > 0$:

$$\mathbb{E} \left[\tilde{\pi}_{n,\alpha}^{\hat{K}} \left(D_\alpha(P_\theta, P^0) > Ms_n | X_1^n \right) \right] \leq \frac{\mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta|X_1^n) \right]}{Ms_n} \leq \frac{r_n}{Ms_n} \xrightarrow{n \rightarrow +\infty} 0.$$

Hence, we obtain the convergence in mean of $\tilde{\pi}_{n,\alpha}^{\hat{K}}(D_\alpha(P_\theta, P^0) > Ms_n | X_1^n)$ to 0, which implies the convergence in probability of $\tilde{\pi}_{n,\alpha}^{\hat{K}}(D_\alpha(P_\theta, P^0) > Ms_n | X_1^n)$ to 0, i.e. the concentration of $\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n)$ to P^0 at rate s_n .

Appendix B. Proof of Theorem 3.

First, we need Donsker and Varadhan's famous variational formula. Refer for example to [Catoni \(2007\)](#) for a proof (Lemma 1.1.3).

Lemma 8 For any probability λ on some measurable space $(\mathbf{E}, \mathcal{E})$ and any measurable function $h : \mathbf{E} \rightarrow \mathbb{R}$ such that $\int e^h d\lambda < \infty$,

$$\log \int e^h d\lambda = \sup_{\rho \in \mathcal{M}_1^+(\mathbf{E})} \left\{ \int h d\rho - \text{KL}(\rho, \lambda) \right\},$$

with the convention $\infty - \infty = -\infty$. Moreover, if h is upper-bounded on the support of λ , then the supremum on the right-hand side is reached by the distribution of the form:

$$\lambda_h(d\beta) = \frac{e^{h(\beta)}}{\int e^h d\lambda} \lambda(d\beta).$$

We come back to the proof of Theorem 3. We adapt the proof of Theorem 4.1 in [Chérif-Abdellatif and Alquier \(2018\)](#).

Proof For any $\alpha \in (0, 1)$ and $\theta \in \Omega := \cup_{K \geq 1} \Theta_K$, using the definition of Rényi divergence and $D_\alpha(P^{\otimes n}, R^{\otimes n}) = nD_\alpha(P, R)$ as data are i.i.d.:

$$\mathbb{E} \left[\exp \left(-\alpha r_n(P_\theta, P^0) + (1 - \alpha)nD_\alpha(P_\theta, P^0) \right) \right] = 1$$

where $r_n(P_\theta, P^0) = \sum_{i=1}^n \log(P^0(X_i)/P_\theta(X_i))$ is the negative log-likelihood ratio. Then we integrate and use Fubini's theorem,

$$\mathbb{E} \left[\int \exp \left(-\alpha r_n(P_\theta, P^0) + (1 - \alpha)nD_\alpha(P_\theta, P^0) \right) \pi(d\theta) \right] = 1.$$

Using Lemma 8,

$$\mathbb{E} \left[\exp \left(\sup_{\rho \in \mathcal{M}_1^+(\Omega)} \left\{ \int \left(-\alpha r_n(P_\theta, P^0) + (1 - \alpha)nD_\alpha(P_\theta, P^0) \right) \rho(d\theta) - \text{KL}(\rho, \pi) \right\} \right) \right] = 1.$$

Then, using Jensen's inequality,

$$\mathbb{E} \left[\sup_{\rho \in \mathcal{M}_1^+(\Omega)} \left\{ \int \left(-\alpha r_n(P_\theta, P^0) + (1 - \alpha)nD_\alpha(P_\theta, P^0) \right) \rho(d\theta) - \text{KL}(\rho, \pi) \right\} \right] \leq 0.$$

Now, we consider $\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n)$ as a distribution on $\mathcal{M}_1^+(\Omega)$ with all its mass on $\Theta_{\hat{K}}$,

$$\mathbb{E} \left[\int \left(-\alpha r_n(P_\theta, P^0) + (1 - \alpha)nD_\alpha(P_\theta, P^0) \right) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta|X_1^n) - \text{KL}(\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n), \pi) \right] \leq 0.$$

We use $\text{KL}(\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n), \pi) = \text{KL}(\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n), \Pi_{\hat{K}}) + \log(\frac{1}{\pi_{\hat{K}}})$, and we rearrange terms:

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta|X_1^n) \right] \\ & \leq \mathbb{E} \left[\frac{\alpha}{1 - \alpha} \int \frac{r_n(P_\theta, P^0)}{n} \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta|X_1^n) + \frac{\text{KL}(\tilde{\pi}_{n,\alpha}^{\hat{K}}(\cdot|X_1^n), \Pi_{\hat{K}})}{n(1 - \alpha)} + \frac{\log(\frac{1}{\pi_{\hat{K}}})}{n(1 - \alpha)} \right]. \end{aligned}$$

By definition of \hat{K} ,

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \\ & \leq \mathbb{E} \left[\inf_{K \geq 1} \left\{ \frac{\alpha}{1-\alpha} \int \frac{r_n(P_\theta, P^0)}{n} \tilde{\pi}_{n,\alpha}^K(d\theta | X_1^n) + \frac{\text{KL}(\tilde{\pi}_{n,\alpha}^K(\cdot | X_1^n), \Pi_K)}{n(1-\alpha)} + \frac{\log(\frac{1}{\pi_K})}{n(1-\alpha)} \right\} \right] \end{aligned}$$

which gives

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \\ & \leq \inf_{K \geq 1} \left\{ \mathbb{E} \left[\frac{\alpha}{1-\alpha} \int \frac{r_n(P_\theta, P^0)}{n} \tilde{\pi}_{n,\alpha}^K(d\theta | X_1^n) + \frac{\text{KL}(\tilde{\pi}_{n,\alpha}^K(\cdot | X_1^n), \Pi_K)}{n(1-\alpha)} + \frac{\log(\frac{1}{\pi_K})}{n(1-\alpha)} \right] \right\} \end{aligned}$$

and hence, by definition of $\tilde{\pi}_{n,\alpha}^K(\cdot | X_1^n)$,

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \\ & \leq \inf_{K \geq 1} \left\{ \mathbb{E} \left[\inf_{\rho \in \mathcal{F}_K} \left\{ \frac{\alpha}{1-\alpha} \int \frac{r_n(P_\theta, P^0)}{n} \rho(d\theta) + \frac{\text{KL}(\rho, \Pi_K)}{n(1-\alpha)} \right\} + \frac{\log(\frac{1}{\pi_K})}{n(1-\alpha)} \right] \right\}. \end{aligned}$$

which leads to,

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \\ & \leq \inf_{K \geq 1} \inf_{\rho \in \mathcal{F}_K} \left\{ \mathbb{E} \left[\frac{\alpha}{1-\alpha} \int \frac{r_n(P_\theta, P^0)}{n} \rho(d\theta) + \frac{\text{KL}(\rho, \Pi_K)}{n(1-\alpha)} + \frac{\log(\frac{1}{\pi_K})}{n(1-\alpha)} \right] \right\}. \end{aligned}$$

Finally,

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \\ & \leq \inf_{K \geq 1} \left\{ \inf_{\rho_K \in \mathcal{F}_K} \left\{ \frac{\alpha}{1-\alpha} \int \text{KL}(P^0, P_{\theta_K}) \rho_K(d\theta_K) + \frac{\text{KL}(\rho_K, \Pi_K)}{n(1-\alpha)} \right\} + \frac{\log(\frac{1}{\pi_K})}{n(1-\alpha)} \right\}. \end{aligned}$$

The theorem is a direct corollary of this inequality as soon as Assumption 3.1 is satisfied. ■

Appendix C. Proof of Theorem 4.

Proof

Fix $\alpha \in (0, 1)$ and let us prove Theorem 4. Let us recall that $\Theta_K(r_{K,n})$ is defined as the set of parameters $\theta_K^* \in \Theta_K$, for which there is a distribution $\rho_{K,n} \in \mathcal{F}_K$ such that:

$$\int \mathbb{E} \left[\log \frac{P_{\theta_K^*}(X_i)}{P_{\theta_K}(X_i)} \right] \rho_{K,n}(d\theta_K) \leq r_{K,n} \quad \text{and} \quad \text{KL}(\rho_{K,n}, \Pi_K) \leq nr_{K,n}.$$

We begin from:

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \\ & \leq \inf_{K \geq 1} \left\{ \inf_{\rho_K \in \mathcal{F}_K} \left\{ \frac{\alpha}{1-\alpha} \int \text{KL}(P^0, P_{\theta_K}) \rho_K(d\theta_K) + \frac{\text{KL}(\rho_K, \Pi_K)}{n(1-\alpha)} \right\} + \frac{\log(\frac{1}{\pi_K})}{n(1-\alpha)} \right\}. \end{aligned}$$

Then, we write for any K , any $\theta_K \in \Theta_K$, $\theta_K^* \in \Theta_K$:

$$\text{KL}(P^0, P_{\theta_K}) = \text{KL}(P^0, P_{\theta_K^*}) + \mathbb{E} \left[\log \frac{P_{\theta_K^*}(X_i)}{P_{\theta_K}(X_i)} \right]$$

which gives:

$$\begin{aligned} & \mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \\ & \leq \inf_{K \geq 1} \left\{ \inf_{\theta_K^* \in \Theta_K} \left\{ \frac{\alpha}{1-\alpha} \text{KL}(P^0, P_{\theta_K^*}) + \inf_{\rho_K \in \mathcal{F}_K} \left\{ \frac{\alpha}{1-\alpha} \int \mathbb{E} \left[\log \frac{P_{\theta_K^*}(X_i)}{P_{\theta_K}(X_i)} \right] \rho_K(d\theta_K) \right. \right. \right. \\ & \qquad \qquad \qquad \left. \left. \left. + \frac{\text{KL}(\rho_K, \Pi_K)}{n(1-\alpha)} \right\} \right\} + \frac{\log(\frac{1}{\pi_K})}{n(1-\alpha)} \right\}. \end{aligned}$$

Hence, using the definition of $\Theta_K(r_{K,n})$ and upper bounding the right-hand-side of the previous inequality by an inf over $\Theta_K(r_{K,n})$, we conclure:

$$\mathbb{E} \left[\int D_\alpha(P_\theta, P^0) \tilde{\pi}_{n,\alpha}^{\hat{K}}(d\theta | X_1^n) \right] \leq \inf_{K \geq 1} \left\{ \frac{\alpha}{1-\alpha} \inf_{\theta^* \in \Theta_K(r_{K,n})} \text{KL}(P^0, P_{\theta_K^*}) + \frac{1+\alpha}{1-\alpha} r_{K,n} + \frac{\log(\frac{1}{\pi_K})}{n(1-\alpha)} \right\}. \quad \blacksquare$$

Appendix D. Proof of Theorem 5.

Proof

We still consider the framework of probabilistic PCA in Section 4. We assume that there exists a true rank K_0 and a matrix $\mathbf{W}_0 \in \mathbb{R}^{d \times K_0}$ with bounded coefficients such that the true distribution of each X_i is $\mathcal{N}(0, \mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d)$, and we place a prior $\Pi_{K_0} = \mathcal{N}(0, s^2 \mathbf{I}_d)^{\otimes K_0}$ and a variational approximation $\rho_{K_0} = \rho^{\otimes K_0}$ on W given $K = K_0$ where we denote $\rho = \mathcal{N}(0, \frac{1}{dn^2} \mathbf{I}_d)$. We recall that $\pi_K = \frac{1}{d}$ for any $K = 1, \dots, d$.

To obtain the rate of convergence $r_n = dK_0 \log(nd)/n$ for probabilistic PCA, we just need to show that the quantities in Assumption 3.1 are upper bounded by r_n (up to a constant) as we have $\log(1/\pi_{K_0})/n$ much smaller than r_n :

$$\int \text{KL} \left(\mathcal{N}(0, \mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d), \mathcal{N}(0, \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d) \right) \rho_{K_0}(d\theta_K) \quad , \quad \frac{\text{KL}(\rho_{K_0}, \Pi_{K_0})}{n}.$$

We have two terms. The first one, i.e. the Kullback-Leibler term, provides a rate of convergence of $dK_0 \log(dn)/n$ as:

$$\begin{aligned} \text{KL}(\rho_{K_0}, \Pi_{K_0}) &= \sum_{j=1}^{K_0} \text{KL}\left(\mathcal{N}\left(0, \frac{1}{dn^2} \mathbf{I}_d\right), \mathcal{N}\left(0, s^2 \mathbf{I}_d\right)\right) \\ &= \frac{K_0}{2} \left(\frac{1}{n^2 s^2} - d + d \log(s^2) + d \log(dn^2) \right) \\ &\leq \frac{K_0}{2n^2 s^2} - \frac{dK_0}{2} + \frac{dK_0 \log(s^2)}{2} + dK_0 \log(dn). \end{aligned}$$

The integral is much more complicated to deal with. We will show that it leads to a rate faster than $dK_0 \log(dn)/n$. If we denote \mathbb{E} the expectation with respect to ρ_{K_0} , then the integral will be equal to:

$$\frac{1}{2} \mathbb{E} \left[\text{Tr} \left((\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} (\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d) \right) \right] - \frac{d}{2} + \frac{1}{2} \mathbb{E} \left[\log \left(\frac{\det(\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d)}{\det(\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d)} \right) \right].$$

The expectation of the log-ratio is easy to upper bound. We denote $\lambda_1, \dots, \lambda_d$ the positive eigenvalues of the positive definite matrix $\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d$. Then for each $j = 1, \dots, d$, $\lambda_j \geq \sigma^2$ and using Jensen's inequality and the log-concavity of the determinant:

$$\begin{aligned} \mathbb{E} \left[\log \left(\det(\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d) \right) \right] &\leq \log \left(\det \left(\mathbb{E}[\mathbf{W} \mathbf{W}^T] + \sigma^2 \mathbf{I}_d \right) \right) \\ &= \log \left(\det \left(\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d + \frac{1}{dn^2} \mathbf{I}_d \right) \right) \\ &= \sum_{j=1}^d \log \left(\lambda_j + \frac{1}{dn^2} \right) \\ &= \sum_{j=1}^d \log(\lambda_j) + \sum_{j=1}^d \log \left(1 + \frac{1}{\lambda_j dn^2} \right) \\ &= \mathbb{E} \left[\log \left(\det(\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d) \right) \right] + \sum_{j=1}^d \log \left(1 + \frac{1}{\lambda_j dn^2} \right) \\ &\leq \mathbb{E} \left[\log \left(\det(\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d) \right) \right] + \sum_{j=1}^d \frac{1}{\lambda_j dn^2} \\ &\leq \mathbb{E} \left[\log \left(\det(\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d) \right) \right] + \frac{1}{n^2 \sigma^2} \end{aligned}$$

and then the expectation of the log-ratio provides a rate of convergence of $1/n^2$:

$$\mathbb{E} \left[\log \left(\frac{\det(\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}_d)}{\det(\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d)} \right) \right] \leq \frac{1}{n^2 \sigma^2}.$$

The remainder can be bounded as follows:

$$\begin{aligned}
 & \mathbb{E} \left[\text{Tr} \left((\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} (\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d) \right) \right] - d \\
 &= \mathbb{E} \left[\text{Tr} \left((\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} (\mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T) \right) \right] \\
 &\leq \mathbb{E} \left[\| (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} \|_F \times \| \mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T \|_F \right] \\
 &\leq \sqrt{d} \mathbb{E} \left[\| (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} \|_2 \times \| \mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T \|_F \right] \\
 &= \sqrt{d} \mathbb{E} \left[\sigma_{\max} ((\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d)^{-1}) \times \| \mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T \|_F \right] \\
 &= \sqrt{d} \mathbb{E} \left[\sigma_{\min} (\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d)^{-1} \times \| \mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T \|_F \right]
 \end{aligned}$$

i.e.

$$\begin{aligned}
 \mathbb{E} \left[\text{Tr} \left((\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} (\mathbf{W}_0 \mathbf{W}_0^T + \sigma^2 \mathbf{I}_d) \right) \right] - d &\leq \sqrt{d} \mathbb{E} \left[(\sigma^2)^{-1} \times \| \mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T \|_F \right] \\
 &= \frac{\sqrt{d}}{\sigma^2} \mathbb{E} \left[\| \mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T \|_F \right]
 \end{aligned}$$

where $\|\cdot\|_F$ is the Frobenius norm on matrices, $\|\cdot\|_2$ the spectral norm, and $\sigma_{\min}(\mathbf{A})$, $\sigma_{\max}(\mathbf{A})$ the lowest and largest singular values of a matrix \mathbf{A} . We use the fact that for a symmetric semi-definite positive matrix: $\sigma_{\max}(\mathbf{A}^{-1}) = (\sigma_{\min}(\mathbf{A}))^{-1}$ and $\sigma_{\min}(\mathbf{A} + \sigma^2 \mathbf{I}_d) \geq \sigma^2$, as well as the inequality $\|\mathbf{A}\|_F \leq \sqrt{d} \|\mathbf{A}\|_2$ for any $d \times d$ matrix \mathbf{A} .

The only thing left to do is to upper bound the expectation of the Frobenius norm of $\mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T$ by a multiple of $\frac{\sqrt{dK_0 \log(dn)}}{n}$. We use the triangle and Cauchy-Schwarz's inequalities:

$$\begin{aligned}
 \mathbb{E} \left[\| \mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W}\mathbf{W}^T \|_F \right] &\leq \mathbb{E} \left[\| \mathbf{W}\mathbf{W}^T - \mathbf{W}\mathbf{W}_0^T \|_F \right] + \mathbb{E} \left[\| \mathbf{W}\mathbf{W}_0^T - \mathbf{W}_0 \mathbf{W}_0^T \|_F \right] \\
 &\leq \mathbb{E} \left[\| \mathbf{W}(\mathbf{W} - \mathbf{W}_0)^T \|_F \right] + \mathbb{E} \left[\| (\mathbf{W} - \mathbf{W}_0) \mathbf{W}_0^T \|_F \right] \\
 &\leq \mathbb{E} \left[\| \mathbf{W} \|_F \| \mathbf{W} - \mathbf{W}_0 \|_F \right] + \mathbb{E} \left[\| \mathbf{W} - \mathbf{W}_0 \|_F \| \mathbf{W}_0 \|_F \right] \\
 &\leq \sqrt{\mathbb{E} [\| \mathbf{W} \|_F^2] \mathbb{E} [\| \mathbf{W} - \mathbf{W}_0 \|_F^2]} + \sqrt{\mathbb{E} [\| \mathbf{W} - \mathbf{W}_0 \|_F^2] \mathbb{E} [\| \mathbf{W}_0 \|_F^2]} \\
 &\leq \sqrt{\mathbb{E} [\| \mathbf{W} \|_F^2] \mathbb{E} [\| \mathbf{W} - \mathbf{W}_0 \|_F^2]} + \| \mathbf{W}_0 \|_F \sqrt{\mathbb{E} [\| \mathbf{W} - \mathbf{W}_0 \|_F^2]}.
 \end{aligned}$$

We can upper bound $\| \mathbf{W}_0 \|_F = \sqrt{\sum_{i=1}^d \sum_{j=1}^{K_0} (\mathbf{W}_0)_{i,j}^2}$ by $\sqrt{dK_0}C$ where C is an upper bound on each of the coefficients of matrix \mathbf{W}_0 .

Also, we can notice that $dn^2 \| \mathbf{W} - \mathbf{W}_0 \|_F^2 = \sum_{i=1}^d \sum_{j=1}^{K_0} (\sqrt{dn}(\mathbf{W}_{i,j} - (\mathbf{W}_0)_{i,j}))^2$ is a sum of squares of independent standard normal random variables. Thus $dn^2 \| \mathbf{W} - \mathbf{W}_0 \|_F^2$

follows a chi-squared distribution with dK_0 degrees of freedom and its expectation is equal to dK_0 . Hence:

$$\mathbb{E}[\|\mathbf{W} - \mathbf{W}_0\|_F^2] = \frac{K_0}{n^2}.$$

Similarly, as $\mathbf{W}_{i,j} - (\mathbf{W}_0)_{i,j}$ is centered, we get:

$$\begin{aligned} \mathbb{E}[\|\mathbf{W}\|_F^2] &= \mathbb{E}\left[\sum_{i=1}^d \sum_{j=1}^{K_0} \mathbf{W}_{i,j}^2\right] \\ &= \sum_{i=1}^d \sum_{j=1}^{K_0} \mathbb{E}\left[(\mathbf{W}_{i,j} - (\mathbf{W}_0)_{i,j})^2 + (\mathbf{W}_0)_{i,j}^2 - 2(\mathbf{W}_0)_{i,j}(\mathbf{W}_{i,j} - (\mathbf{W}_0)_{i,j})\right] \\ &= \mathbb{E}[\|\mathbf{W} - \mathbf{W}_0\|_F^2] + \|\mathbf{W}_0\|_F^2 \\ &\leq \frac{K_0}{n^2} + dK_0C^2 \\ &= \left(dC^2 + \frac{1}{n^2}\right)K_0. \end{aligned}$$

Thus, we obtain:

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W} \mathbf{W}^T\|_F\right] &\leq \frac{\sqrt{K_0}}{n} \sqrt{K_0} \sqrt{dC^2 + \frac{1}{n^2}} + \sqrt{dK_0}C \frac{\sqrt{K_0}}{n} \\ &= \frac{K_0}{n} \sqrt{dC^2 + \frac{1}{n^2}} + \frac{\sqrt{dK_0}C}{n} \\ &\leq \frac{K_0}{n} \left(\sqrt{d}C + \frac{1}{n}\right) + \frac{\sqrt{dK_0}C}{n} \\ &= \frac{K_0}{n} \left(2\sqrt{d}C + \frac{1}{n}\right). \end{aligned}$$

Hence, the order of the upper bound of the expectation of the Fobrenius norm of matrix $\mathbf{W}_0 \mathbf{W}_0^T - \mathbf{W} \mathbf{W}^T$ is $\frac{\sqrt{dK_0}}{n} < \frac{\sqrt{dK_0} \log(dn)}{n}$.

Finally, the consistency rate associated with the integral term is $\frac{dK_0}{n}$, and the overall rate of convergence is $\frac{dK_0 \log(dn)}{n}$. \blacksquare

Appendix E. Computation of the ELBO for probabilistic PCA.

We consider the framework of probabilistic PCA detailed in Section 4. Given rank K , we place independent Gaussian priors on the columns W_1, \dots, W_K of \mathbf{W} such that $\Pi_K = \mathcal{N}(0, s^2 \mathbf{I}_d)^{\otimes K}$, and Gaussian independent variational approximations $\mathcal{N}(\mu_j, \Sigma_j)$ for the columns of \mathbf{W} . The ELBO associated with rank K and variational approximation $\rho_K = \otimes_{j=1}^K \mathcal{N}(\mu_j, \Sigma_j)$ is given by:

$$\text{ELBO}_K(\rho_K) = \alpha \int \ell_n(\mathbf{W}) \rho_K(d\mathbf{W}) - \text{KL}(\rho_K, \Pi_K).$$

The Kullback-Leibler term $\text{KL}(\rho_K, \Pi_K)$ is equal to:

$$\frac{1}{2} \sum_{j=1}^K \left\{ \frac{\text{Tr}(\boldsymbol{\Sigma}_j)}{s^2} + \frac{\mu_j^T \mu_j}{s^2} - \log(\det(\boldsymbol{\Sigma}_j)) \right\} - \frac{dK}{2} + \frac{dK \log(s^2)}{2}$$

while the average log-likelihood $\int \ell_n(\mathbf{W}) \rho_K(d\mathbf{W})$ is:

$$-\frac{dn}{2} \log(2\pi) - \frac{n}{2} \int \log(\det(\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)) \rho_K(d\mathbf{W}) - \frac{1}{2} \sum_{i=1}^n \int X_i^T (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} X_i \rho_K(d\mathbf{W})$$

where both integrals can be computed thanks to Monte-Carlo sampling approximations:

$$\int \log(\det(\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)) \rho_K(d\mathbf{W}) \approx \sum_{\ell=1}^N \log(\det(\mathbf{W}^{(\ell)} \mathbf{W}^{(\ell)T} + \sigma^2 \mathbf{I}_d))$$

and

$$\int X_i^T (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} X_i \rho_K(d\mathbf{W}) \approx \sum_{\ell=1}^N X_i^T (\mathbf{W}^{(\ell)} \mathbf{W}^{(\ell)T} + \sigma^2 \mathbf{I}_d)^{-1} X_i$$

where $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(N)}$ are N i.i.d. data sampled from ρ_K .

The inverse matrix $(\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1}$ can be derived thanks to classical inversion algorithms. For instance, it is possible to do so in $\mathcal{O}(Kd^2)$ operations instead of the classical $\mathcal{O}(d^3)$ inversion procedure thanks to Sherman-Morrison formula: for any matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ and vectors $u, v \in \mathbb{R}^d$ such that $\mathbf{A} + uv^T$ is invertible,

$$(\mathbf{A} + uv^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} uv^T \mathbf{A}^{-1}}{1 + v^T \mathbf{A}^{-1} u}.$$

We write

$$\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d = \sigma^2 \mathbf{I}_d + \sum_{j=1}^K W_j W_j^T = \left(\sigma^2 \mathbf{I}_d + \sum_{j=1}^{K-1} W_j W_j^T \right) + W_K W_K^T$$

and iterate K times Sherman-Morrison formula. The first time, we apply it to $\mathbf{A} = \sigma^2 \mathbf{I}_d + \sum_{j=1}^{K-1} W_j W_j^T$ and $u = v = \mathbf{W}_K$, then to $\mathbf{A} = \sigma^2 \mathbf{I}_d + \sum_{j=1}^{K-2} W_j W_j^T$ and $u = v = \mathbf{W}_{K-1}$, and so on. We finally obtain $(\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d)^{-1} = \mathbf{M}_K$ where:

$$\begin{cases} \mathbf{M}_0 = \sigma^2 \mathbf{I}_d \\ \forall j = 1, \dots, K, \mathbf{M}_j = \mathbf{M}_{j-1} - \frac{1}{1 + \mathbf{W}_{j-1}^T \mathbf{M}_{j-1} \mathbf{W}_{j-1}} \mathbf{W}_{j-1} \mathbf{W}_{j-1}^T \text{ with } \mathbf{Z}_j = \mathbf{M}_{j-1} \mathbf{W}_j. \end{cases}$$

In order to compute the maximum value $\text{ELBO}(K)$ of the ELBO associated with rank K , one can use a stochastic gradient descent on $(\mu_1, \boldsymbol{\Sigma}_1, \dots, \mu_K, \boldsymbol{\Sigma}_K)$ that will converge to a local maximum and will give the variational estimator for rank K . Then, maximizing $\text{ELBO}(K)$ over desired values of K leads to the optimal number of principal components and to the associated optimal variational approximation.

Appendix F. Results in matrix norm for probabilistic PCA.

To prove Corollaries 6 and 7, we need the two lemmas presented behind. We introduce some notations first. We refer the interested reader to Forth et al. (2014) for more details.

Notations : Let us call \mathcal{S}_d^+ the set of $d \times d$ symmetric positive semi-definite matrices, and $\mathcal{X}_M = \{\mathbf{A} \in \mathcal{S}_d^+ / \|\mathbf{A}\|_2 \leq M\}$. We define the vectorization of Matrix $\mathbf{A} \in \mathbb{R}^{p \times q}$ with columns X_1, \dots, X_q :

$$\text{Vec}(\mathbf{A}) = (\mathbf{A}_1^T, \dots, \mathbf{A}_q^T)^T \in \mathbb{R}^{p \times q}.$$

We define the Frobenius inner product of two matrices $\mathbf{A} \in \mathbb{R}^{p \times q}$ and $\tilde{\mathbf{A}} \in \mathbb{R}^{p \times q}$, that is the sum of componentwise products:

$$\mathbf{A} \cdot \tilde{\mathbf{A}} = \text{Vec}(\mathbf{A})^T \text{Vec}(\tilde{\mathbf{A}}).$$

Notice that $\|\mathbf{A}\|_F^2 = \mathbf{A} \cdot \mathbf{A} = \text{Vec}(\mathbf{A})^T \text{Vec}(\mathbf{A})$.

We also introduce the Kronecker and Box products of two matrices $\mathbf{A} \in \mathbb{R}^{p_1 \times q_1}$ and $\tilde{\mathbf{A}} \in \mathbb{R}^{p_2 \times q_2}$ which are respectively the matrices $\mathbf{A} \otimes \tilde{\mathbf{A}} \in \mathbb{R}^{p_1 p_2 \times q_1 q_2}$ and $\mathbf{A} \boxtimes \tilde{\mathbf{A}} \in \mathbb{R}^{p_1 p_2 \times q_1 q_2}$ such that their coefficients are defined as:

$$\begin{aligned} (\mathbf{A} \otimes \tilde{\mathbf{A}})_{p_2(i-1)+j, q_2(k-1)+l} &= \mathbf{A}_{i,k} \tilde{\mathbf{A}}_{j,l}, \\ (\mathbf{A} \boxtimes \tilde{\mathbf{A}})_{p_2(i-1)+j, q_1(k-1)+l} &= \mathbf{A}_{i,l} \tilde{\mathbf{A}}_{j,k} \end{aligned}$$

for any integers i, j, k, l such that $1 \leq i \leq p_1$, $1 \leq j \leq q_1$, $1 \leq k \leq p_2$, $1 \leq l \leq q_2$.

We have the following properties for any matrix \mathbf{P} :

$$\begin{aligned} (\mathbf{A} \otimes \tilde{\mathbf{A}}) \text{Vec}(\mathbf{P}) &= \text{Vec}(\tilde{\mathbf{A}} \mathbf{P} \mathbf{A}^T), \\ (\mathbf{A} \boxtimes \tilde{\mathbf{A}}) \text{Vec}(\mathbf{P}) &= \text{Vec}(\tilde{\mathbf{A}} \mathbf{P}^T \mathbf{A}^T). \end{aligned}$$

We also define the gradient $\nabla f(\mathbf{A}) \in \mathbb{R}^{p \times q}$ and the Hessian $\nabla^2 f(\mathbf{A}) \in \mathbb{R}^{pq \times pq}$ of a differentiable function $f : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}$ at matrix \mathbf{A} :

$$\begin{aligned} (\nabla f(\mathbf{A}))_{p_2(i-1)+j, q_2(k-1)+l} &= \frac{\partial f(\mathbf{A})}{\partial \mathbf{A}_{i,j}}, \\ (\nabla^2 f(\mathbf{A}))_{p_2(j-1)+i, p_2(l-1)+k} &= \frac{\partial^2 f(\mathbf{A})}{\partial \mathbf{A}_{i,j} \partial \mathbf{A}_{k,l}} \end{aligned}$$

for any integers i, j, k, l such that $1 \leq i, k \leq p$, $1 \leq j, l \leq q$ where ∂f is the partial derivative of f .

We say that a differentiable function $f : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}$ is s -strongly convex in $\mathcal{S} \subset \mathbb{R}^{pq \times pq}$ with respect to the norm $\|\cdot\|$ as soon as one of the two following equivalent properties is satisfied:

$$f(\mathbf{A}) \geq f(\tilde{\mathbf{A}}) + \nabla f(\tilde{\mathbf{A}}) \cdot (\mathbf{A} - \tilde{\mathbf{A}}) + \frac{s}{2} \|\mathbf{A} - \tilde{\mathbf{A}}\|^2$$

or

$$\text{Vec}(\mathbf{P})^T \nabla^2 f(\mathbf{A}) \text{Vec}(\mathbf{P}) \geq s \|\mathbf{P}\|^2$$

for any matrix $\mathbf{A}, \tilde{\mathbf{A}} \in \mathcal{S}$ and any symmetric matrix $\mathbf{P} \in \mathbb{R}^{pq \times pq}$.

Lemma 9 *Then, function $f : \mathbf{A} \rightarrow -\log(\det(\mathbf{A} + M\mathbf{I}_d))$ is $1/(M + \sigma^2)^2$ strongly convex in \mathcal{X}_M with respect to the Frobenius norm.*

Proof The proof follows the same steps than the proof of Theorem 3.1 in [Moridomi et al. \(2018\)](#).

The Hessian of function f at any symmetric matrix in $\mathbf{A} \in \mathcal{X}_M$ is given by (see [Forth et al. \(2014\)](#)):

$$\nabla^2 f(\mathbf{A}) = \left((\mathbf{A} + M\mathbf{I}_d)^{-1} \right)^T \boxtimes (\mathbf{A} + M\mathbf{I}_d)^{-1} = (\mathbf{A} + M\mathbf{I}_d)^{-1} \boxtimes (\mathbf{A} + M\mathbf{I}_d)^{-1}.$$

Then, we have for any $\mathbf{A} \in \mathcal{X}_M$ and any symmetric matrix $\mathbf{P} \in \mathbb{R}^{pq \times pq}$:

$$\begin{aligned} \text{Vec}(\mathbf{P})^T \nabla^2 f(\mathbf{A}) \text{Vec}(\mathbf{P}) &= \text{Vec}(\mathbf{P})^T \left((\mathbf{A} + M\mathbf{I}_d)^{-1} \boxtimes (\mathbf{A} + M\mathbf{I}_d)^{-1} \right) \text{Vec}(\mathbf{P}) \\ &= \text{Vec}(\mathbf{P})^T \text{Vec} \left((\mathbf{A} + M\mathbf{I}_d)^{-1} \mathbf{P}^T (\mathbf{A} + M\mathbf{I}_d)^{-1} \right) \\ &= \text{Vec}(\mathbf{P})^T \text{Vec} \left((\mathbf{A} + M\mathbf{I}_d)^{-1} \mathbf{P} (\mathbf{A} + M\mathbf{I}_d)^{-1} \right) \\ &= \text{Vec}(\mathbf{P})^T \left((\mathbf{A} + M\mathbf{I}_d)^{-1} \otimes (\mathbf{A} + M\mathbf{I}_d)^{-1} \right) \text{Vec}(\mathbf{P}). \end{aligned}$$

Note that the eigenvalues of a Kronecker product $\mathbf{A} \otimes \mathbf{P}$ are the products of an eigenvalue of \mathbf{A} and an eigenvalue of \mathbf{P} , and the eigenvalues of \mathbf{P}^{-1} are the inverse of the eigenvalues of \mathbf{P} . Moreover, the maximum eigenvalue of $\mathbf{A} + M\mathbf{I}_d$ is $\|\mathbf{A}\|_2 + \sigma^2$, so the minimum eigenvalue of $(\mathbf{A} + M\mathbf{I}_d)^{-1} \otimes (\mathbf{A} + M\mathbf{I}_d)^{-1}$ is equal to $(\|\mathbf{A}\|_2 + \sigma^2)^{-2}$. Hence, for any matrix $\mathbf{A} \in \mathcal{X}_M$, we get:

$$\begin{aligned} \text{Vec}(\mathbf{P})^T \left((\mathbf{A} + M\mathbf{I}_d)^{-1} \otimes (\mathbf{A} + M\mathbf{I}_d)^{-1} \right) \text{Vec}(\mathbf{P}) &\geq (\|\mathbf{A}\|_2 + \sigma^2)^{-2} \text{Vec}(\mathbf{P})^T \text{Vec}(\mathbf{P}) \\ &\geq \frac{1}{(M + \sigma^2)^2} \text{Vec}(\mathbf{P})^T \text{Vec}(\mathbf{P}), \end{aligned}$$

and we conclude using the definition of the strong convexity and $\|\mathbf{P}\|_F^2 = \text{Vec}(\mathbf{P})^T \text{Vec}(\mathbf{P})$. \blacksquare

Lemma 10 For any $\alpha \in (0, 1)$ and any matrices $\mathbf{W} \in \mathbb{R}^{d \times K_1}$ and $\tilde{\mathbf{W}} \in \mathbb{R}^{d \times K_2}$, as soon as the spectral norms of $\mathbf{W}\mathbf{W}^T$ and $\tilde{\mathbf{W}}\tilde{\mathbf{W}}^T$ are bounded by a constant B^2 , then:

$$D_\alpha(P_{\mathbf{W}}, P_{\tilde{\mathbf{W}}}) \geq \frac{\alpha}{16(B^2 + \sigma^2)^2} \left\| \tilde{\mathbf{W}}\tilde{\mathbf{W}}^T - \mathbf{W}\mathbf{W}^T \right\|_F^2.$$

Proof

We recall that function $f : \mathbf{A} \rightarrow -\log(\det(\mathbf{A} + M\mathbf{I}_d))$ is $1/(M + \sigma^2)^2$ strongly convex in \mathcal{X}_M with respect to the Fobrenius norm according to Lemma 9. Hence, for any matrices \mathbf{A} and $\tilde{\mathbf{A}}$ in \mathcal{X}_M , we have:

$$\begin{aligned} -\log(\det((1 - \alpha)\mathbf{A} + \alpha\tilde{\mathbf{A}})) &\leq -(1 - \alpha)\log(\det(\mathbf{A})) - \alpha\log(\det(\tilde{\mathbf{A}})) \\ &\quad - \frac{1}{2}\alpha(1 - \alpha)\frac{1}{4M^2}\|\tilde{\mathbf{A}} - \mathbf{A}\|_F^2. \end{aligned}$$

We rearrange terms:

$$\log \left(\frac{\det((1-\alpha)\mathbf{A} + \alpha\tilde{\mathbf{A}})}{\det(\mathbf{A})^{1-\alpha} \det(\tilde{\mathbf{A}})^\alpha} \right) \geq \frac{\alpha(1-\alpha)}{8M^2} \|\tilde{\mathbf{A}} - \mathbf{A}\|_F^2.$$

Now, we use the fact that:

$$D_\alpha(\mathcal{N}(0, \mathbf{A}), \mathcal{N}(0, \tilde{\mathbf{A}})) = \frac{1}{2(1-\alpha)} \log \left(\frac{\det((1-\alpha)\mathbf{A} + \alpha\tilde{\mathbf{A}})}{\det(\mathbf{A})^{1-\alpha} \det(\tilde{\mathbf{A}})^\alpha} \right)$$

to get for any matrices $\mathbf{W} \in \mathbb{R}^{d \times K_1}$ and $\tilde{\mathbf{W}} \in \mathbb{R}^{d \times K_2}$ such that $\|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^T + \sigma^2 \mathbf{I}_d\|_2 \leq M$ and $\|\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d\|_F \leq M$:

$$D_\alpha(P_{\mathbf{W}}, P_{\tilde{\mathbf{W}}}) \geq \frac{\alpha}{16M^2} \|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^T + \sigma^2 \mathbf{I}_d - \mathbf{W}\mathbf{W}^T - \sigma^2 \mathbf{I}_d\|_F^2.$$

Moreover, for any matrix $\mathbf{W} \in \mathbb{R}^{d \times K}$ such that the spectral norm of $\mathbf{W}\mathbf{W}^T$ is bounded by B^2 , we have $\|\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d\|_2 \leq B^2 + \sigma^2$. We conclude using the previous inequality for $M = B^2 + \sigma^2$. ■

Now, let us go back to the proof of Corollary 6.

Proof

We assume that there exists a true model \mathcal{M}_{K_0} such that $P^0 = P_{\mathbf{W}_0}$ with $\mathbf{W}_0 \in \mathbb{R}^{d \times K_0}$ and such that the spectral norm of \mathbf{W}_0 is bounded by B (hence the coefficients of \mathbf{W}_0 are also bounded). As clip_B is a projection onto a closed convex set with respect to the Frobenius norm, we have for any matrix $\mathbf{W} \in \mathbb{R}^{d \times \hat{K}}$:

$$\|\text{clip}_B(\mathbf{W}\mathbf{W}^T) - \text{clip}_B(\mathbf{W}_0\mathbf{W}_0^T)\|_F \leq \|\mathbf{W}\mathbf{W}^T - \mathbf{W}_0\mathbf{W}_0^T\|_F$$

and as the coefficients of $\mathbf{W}_0\mathbf{W}_0^T$ are bounded by B^2 :

$$\|\text{clip}_B(\mathbf{W}\mathbf{W}^T) - \mathbf{W}_0\mathbf{W}_0^T\|_F = \|\text{clip}_B(\mathbf{W}\mathbf{W}^T) - \text{clip}_B(\mathbf{W}_0\mathbf{W}_0^T)\|_F.$$

According to Lemma 10, we get for any matrix $\mathbf{W} \in \mathbb{R}^{d \times \hat{K}}$:

$$\|\text{clip}_B(\mathbf{W}\mathbf{W}^T) - \mathbf{W}_0\mathbf{W}_0^T\|_F^2 \leq \frac{16(B^2 + \sigma^2)^2}{\alpha} D_\alpha(P_{\mathbf{W}}, P_{\mathbf{W}_0}).$$

Thus:

$$\begin{aligned} \mathbb{E} \left[\int \|\text{clip}_B(\mathbf{W}\mathbf{W}^T) - \mathbf{W}_0\mathbf{W}_0^T\|_F^2 \tilde{\pi}_{n,\alpha}^{\hat{K}}(dW|X_1^n) \right] \\ \leq \frac{16(B^2 + \sigma^2)^2}{\alpha} \mathbb{E} \left[\int D_\alpha(P_{\mathbf{W}}, P_{\mathbf{W}_0}) \tilde{\pi}_{n,\alpha}^{\hat{K}}(dW|X_1^n) \right] \end{aligned}$$

and we use Theorem 5:

$$\mathbb{E} \left[\int D_\alpha(P_{\mathbf{W}}, P_{\mathbf{W}_0}) \tilde{\pi}_{n,\alpha}^{\hat{K}}(dW|X_1^n) \right] = \mathcal{O} \left(\frac{dK_0 \log(dn)}{n} \right).$$

which ends the proof. ■

We can obtain Corollary 7 using a simple convexity argument.