

# Exact Dimensionality Selection for Bayesian PCA

Charles Bouveyron, Pierre Latouche, Pierre-Alexandre Mattei

► **To cite this version:**

Charles Bouveyron, Pierre Latouche, Pierre-Alexandre Mattei. Exact Dimensionality Selection for Bayesian PCA. 2017. <hal-01484099>

**HAL Id: hal-01484099**

**<https://hal.archives-ouvertes.fr/hal-01484099>**

Submitted on 6 Mar 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## EXACT DIMENSIONALITY SELECTION FOR BAYESIAN PCA

**Charles Bouveyron***Laboratoire MAP5, UMR CNRS 8145**Université Paris Descartes & Sorbonne Paris Cité**charles.bouveyron@parisdescartes.fr***Pierre Latouche***Laboratoire SAMM, EA 4543**Université Paris 1 Panthéon-Sorbonne**pierre.latouche@univ-paris1.fr***Pierre-Alexandre Mattei***Laboratoire MAP5, UMR CNRS 8145**Université Paris Descartes & Sorbonne Paris Cité**pierre-alexandre.mattei@parisdescartes.fr***Abstract**

We present a Bayesian model selection approach to estimate the intrinsic dimensionality of a high-dimensional dataset. To this end, we introduce a novel formulation of the probabilistic principal component analysis model based on a normal-gamma prior distribution. In this context, we exhibit a closed-form expression of the marginal likelihood which allows to infer an optimal number of components. We also propose a heuristic based on the expected shape of the marginal likelihood curve in order to choose the hyperparameters. In non-asymptotic frameworks, we show on simulated data that this exact dimensionality selection approach is competitive with both Bayesian and frequentist state-of-the-art methods.

**Keywords:** Dimensionality reduction, Marginal likelihood, Multivariate analysis, Model selection, Principal components.

**1. Introduction**

The computer age is characterized by a surge of multivariate data, which is often difficult to explore or describe. A natural way to deal with such datasets is to reduce their dimensionality in a interpretable way, trying not to lose too much information. Accordingly, a wide range of dimension-reduction techniques have been developed over the years. Principal component analysis (PCA), perhaps the earliest of these techniques, remains today one of the most widely used (Jolliffe and Cadima, 2016). Introduced by Pearson (1901) and rediscovered by Hotelling (1933) in the beginning of the twentieth century, PCA has had indeed an ubiquitous role in statistical analysis since the introduction of electronic computation in the 1950s. Recent examples include climate research (Hannachi et al., 2006), genome-wide expression studies (Ringnér, 2008), massive text mining (Zhang and El Ghaoui, 2011) and even deep learning (Chan et al., 2015). For a more exhaustive overview of past applications of PCA, we defer the reader to the monograph of Jolliffe (2002) or the recent review paper of Jolliffe and Cadima (2016).

Specifically, PCA consists in a simple procedure: the practitioner orthogonally projects his multivariate data on a space spanned by the eigenvectors associated with the largest eigenvalues of the empirical covariance matrix. The dimension of the representation learnt in this way is simply the number of eigenvectors – called principal components (PCs) – kept for the projection. However, it may come as a surprise that in spite of the popularity of this method, no authoritative solution has been widely accepted for choosing how many

PCs should be computed. Common practice is to choose this dimension by considering the eigenvalues scree of the sample covariance matrix. This ad-hoc technique, popularized by Cattell (1966), has been largely modified and perfected over the last fifty years (Jackson, 1993; Zhu and Ghodsi, 2006), and is often chosen when PCA is used as a building block within a larger algorithmic framework – see e.g. Bouveyron et al. (2007) for an example in cluster analysis or Evangelopoulos et al. (2012) in latent semantic analysis. However, more refined approaches have also been developed. Earlier works were based on hypothesis testing (Jolliffe, 2002, section 6.1.4). Cross-validation, suggested by Wold (1978) and developed over the years (Bro et al., 2008), is known to be effective in a wide variety of settings (Josse and Husson, 2012). Another fruitful line of work follows the seminal article of Tipping and Bishop (1999), who recast PCA as a simple inferential problem. Their model, called probabilistic PCA (PPCA), led to several model-based methods for dimensionality selection, both from frequentist (Ulfarsson and Solo, 2008; Bouveyron et al., 2011; Passemier et al., 2017) and Bayesian (Bishop, 1999; Minka, 2000; Hoyle, 2008) perspectives.

Most of the aforementioned methods are based on asymptotic considerations. However, it was recently proven that, in an asymptotic framework, hard thresholding the eigenvalues surprisingly suffices to provide an optimal dimensionality (Gavish and Donoho, 2014). Thus, the path to more efficient schemes for finding the number of PCs goes through the study of non-asymptotic criteria, which have been overlooked in the past. A natural non-asymptotic answer is provided by exact Bayesian model selection, which was previously used at the price of computationally expensive Markov chain Monte Carlo (MCMC) sampling (Hoff, 2007). We present here a prior structure based on the PPCA model that allows us to exhibit a closed-form expression of the marginal likelihood, leading to an efficient algorithm that selects the number of PCs without any asymptotic assumption. Specifically, we rely on a normal prior distribution over the loading matrix and a gamma prior distribution over the noise variance. Imposing a simple constraint on the hyperparameters of the respective distributions, we show that this allows the data to marginally follow a generalized Laplace distribution, leading to an efficient closed-form computation of the marginal likelihood. We also propose a heuristic based on the expected shape of the marginal likelihood curve in order to choose hyperparameters. With simulated data, we demonstrate that our approach is competitive compared to state-of-the-art methods, especially in non asymptotic settings and with less observations than variables. This setting is at the core of many practical problems, such as genomics and chemometrics.

In Section 2, we briefly review PPCA and present several dimensionality selection techniques based on this model. The new normal-gamma prior is presented in Section 3 together with a derivation of the closed-form expression of the marginal likelihood. A heuristic to choose hyperparameters is also presented. Numerical experiments are provided in Section 4.

## 2. Choosing the intrinsic dimension in probabilistic PCA

Let us assume that a centered independent and identically distributed (i.i.d.) sample  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  is observed that we aim at projecting onto a  $d$ -dimensional subspace while retaining as much variance as possible. All the observations are stored in the  $n \times p$  matrix  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ .

## 2.1 Probabilistic PCA

The PPCA model  $\mathcal{M}_d$  assumes that, for all  $i \in \{1, \dots, n\}$ , each observation is driven by the following generative model

$$\mathbf{x}_i = \mathbf{W}\mathbf{y}_i + \boldsymbol{\varepsilon}_i, \quad (1)$$

where  $\mathbf{y}_i \sim \mathcal{N}(0, \mathbf{I}_d)$  is a low-dimensional Gaussian latent vector,  $\mathbf{W}$  is a  $p \times d$  parameter matrix called the *loading matrix* and  $\boldsymbol{\varepsilon}_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_p)$  is a Gaussian noise term.

This model is an instance of factor analysis and was first introduced by Lawley (1953). Tipping and Bishop (1999) then presented a thorough study of this model. In particular, expanding a result of Theobald (1975), they proved that this generative model is indeed equivalent to PCA in the sense that the principal components of  $\mathbf{X}$  can be retrieved using the maximum likelihood (ML) estimator  $\mathbf{W}_{\text{ML}}$  of  $\mathbf{W}$ . More specifically, if  $\mathbf{A}$  is the  $p \times d$  matrix of ordered principal eigenvectors of  $\mathbf{X}^T \mathbf{X}$  and if  $\boldsymbol{\Lambda}$  is the  $d \times d$  diagonal matrix with corresponding eigenvalues, we have

$$\mathbf{W}_{\text{ML}} = \mathbf{A}(\boldsymbol{\Lambda} - \sigma^2 \mathbf{I}_d)^{1/2} \mathbf{R}, \quad (2)$$

where  $\mathbf{R}$  is an arbitrary orthogonal matrix.

Under this sound probabilistic framework, dimension selection can be recast as a *model selection problem*, for which standard techniques are available. We review a few important ones in the next subsection.

## 2.2 Model selection for PPCA

The problem of finding an appropriate dimension can be seen as choosing a "best model" within a family of models  $(\mathcal{M}_d)_{d \in \{1, \dots, p-1\}}$ . A first popular approach would be to use likelihood penalization, leading to the choice

$$d^* \in \operatorname{argmax}_{d \in \{1, \dots, p-1\}} \{\log p(\mathbf{X} | \mathbf{W}_{\text{ML}}, \sigma_{\text{ML}}, \mathcal{M}_d) - \operatorname{pen}(d)\},$$

where  $\operatorname{pen}$  is a penalty which grows with  $d$ . These methods include the popular Akaike information criterion (AIC, Akaike, 1974), the Bayesian information criterion (BIC, Schwarz, 1978), or other refined approaches (Bai and Ng, 2002). However, their merits are mainly asymptotic, and our main interest in this paper is to investigate non-asymptotic scenarios. While the penalty term is usually necessary to avoid selecting the largest model, under a constrained PPCA model, called isotropic PPCA, Bouveyron et al. (2011) proved that regular maximum likelihood was suprisingly consistent. While the theoretical optimality of this method is also asymptotic, the fact that it directly maximizes a likelihood criterion which is not derived based on asymptotic considerations makes it of particular interest within the scope of this paper.

Another interesting set of techniques non-asymptotic in essence is Bayesian model selection (Kass and Raftery, 1995). While BIC does not actually approximate the marginal likelihood in the case of PPCA because of violated regularity conditions (Drton and Plummer, 2017), a more refined approach was followed by Minka (2000) who derived a Laplace approximation of the marginal likelihood. This technique, albeit asymptotic, has been proven empirically efficient in several small-sample scenarios.

Another interesting framework considered in the literature is the case where both  $n$  and  $p$  grow to infinity. Several consistent estimators have been proposed, both from a penalization point of view (Bai and Ng, 2002; Passemier et al., 2017), using Stein’s unbiased risk estimator (Ulfarsson and Solo, 2008) or in a Bayesian context (Hoyle, 2008). While these high-dimensional scenarios are of growing importance, they fall beyond the scope of this paper, which is focused on the non-asymptotic setting (with potentially fewer observations than variables), for which very few automatic dimension selection methods are available.

### 3. Exact dimensionality selection for PPCA under a normal-gamma prior

In this section, we present a prior structure that leads to a closed-form expression for the marginal likelihood of PPCA.

#### 3.1 The model

We consider the regular PPCA model already defined in (1),

$$\forall i \in \{1, \dots, n\}, \mathbf{x}_i = \mathbf{W}\mathbf{y}_i + \boldsymbol{\varepsilon}_i,$$

where  $\mathbf{y}_i \sim \mathcal{N}(0, \mathbf{I}_d)$ ,  $\mathbf{W}$  is a  $p \times d$  parameter matrix, and  $\boldsymbol{\varepsilon}_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_p)$ . We rely on a Gaussian prior distribution over the loading matrix  $\mathbf{W}$  and a gamma prior distribution over the noise variance  $\sigma^2$ . Specifically, we use a gamma prior  $\sigma^2 \sim \text{Gamma}(a, b)$  with hyperparameters  $a > 0$  and  $b > 0$  together with i.i.d. Gaussian priors  $w_{jk} \sim \mathcal{N}(0, \phi^{-1})$  for  $j \in \{1, \dots, p\}$  and  $k \in \{1, \dots, d\}$  with some precision hyperparameter  $\phi > 0$ .

Within the framework of Bayesian model uncertainty (Kass and Raftery, 1995), the posterior probabilities of models can be written as, for all  $d \in \{1, \dots, p\}$ ,

$$p(\mathcal{M}_d | \mathbf{X}, a, b, \phi) \propto p(\mathbf{X} | a, b, \phi, \mathcal{M}_d) p(\mathcal{M}_d), \quad (3)$$

where

$$p(\mathbf{X} | a, b, \phi, \mathcal{M}_d) = \int_{\mathbb{R}^{d \times p} \times \mathbb{R}^+} p(\mathbf{X} | \mathbf{W}, \sigma, \mathcal{M}_d) p(\mathbf{W} | \phi) p(\sigma | a, b) d\mathbf{W} d\sigma,$$

is the *marginal likelihood* of the data. Note that this expression also involves model prior probabilities – in this paper, we will simply consider a uniform prior

$$\forall d \in \{1, \dots, p\}, p(\mathcal{M}_d) \propto 1.$$

Computing the high-dimensional integral of the marginal likelihood usually comes at the price of various approximations (Bishop, 1999; Minka, 2000; Hoyle, 2008) or expensive sampling (Hoff, 2007). However, with our specific choice of priors, and imposing a constraint on their respective hyperparameters, we obtain a closed-form expression for the marginal likelihood.

**Theorem 1** *Let  $d \in \{1, \dots, p\}$ . Under the normal-gamma prior with  $b = \phi/2$ , the log-marginal likelihood of model  $\mathcal{M}_d$  is given by*

$$\begin{aligned} \log p(\mathbf{X}|a, \phi, \mathcal{M}_d) &= \sum_{i=1}^n \log p(\mathbf{x}_i|a, \phi, \mathcal{M}_d) \\ &= -\frac{np}{2} \log(2\pi) - \frac{np}{2} \log(2\phi^{-1}) - n \log \Gamma(a + d/2) \\ &\quad + (a + \frac{d-p}{2}) \sum_{i=1}^n \log\left(\frac{\sqrt{\phi} \|\mathbf{x}_i\|_2}{2}\right) + \sum_{i=1}^n \log K_{a+(d-p)/2}(\sqrt{\phi} \|\mathbf{x}_i\|_2), \end{aligned} \quad (4)$$

where  $K_\nu$  is the modified Bessel function of the second kind of order  $\nu \in \mathbb{R}$ .

A detailed proof of this theorem is given in the next subsection.

To the best of our knowledge, this result is the first computation of the marginal likelihood of a PPCA model. It is worth mentioning that, in a slightly different context, Ando (2009) also derived the marginal likelihood of a factor analysis model, with Student factors. Similarly, Bouveyron et al. (2016) derived the exact marginal likelihood of the noiseless PPCA model, in order to obtain sparse PCs.

While Gaussian priors for the loading matrix have been extensively used in the past (Bishop, 1999; Archambeau and Bach, 2009; Bouveyron et al., 2016), it is worth noticing that the use of a gamma prior for a variance parameter is rather peculiar. Indeed, most Bayesian hierarchical models choose *inverse-gamma* priors for variances. This choice is often motivated by its conjugacy properties (see e.g. George and McCulloch, 1993, for a linear regression example or Murphy, 2007, in a wider setting). The derivation provided in the next subsection notably explains why this gamma prior over  $\sigma^2$  actually arises naturally.

### 3.2 Derivation of the marginal likelihood

We begin by shortly reviewing the generalized Laplace distribution, which will prove to be key within the PPCA framework. This distribution was introduced by Kotz et al. (2001, p. 257). For a more detailed overview, see Kozubowski et al. (2013).

**Definition 2** *A random variable  $\mathbf{z} \in \mathbb{R}^p$  is said to have a **multivariate generalized asymmetric Laplace distribution** with parameters  $s > 0$ ,  $\boldsymbol{\mu} \in \mathbb{R}^p$  and  $\boldsymbol{\Sigma} \in \mathcal{S}_p^+$  if its characteristic function is*

$$\forall \mathbf{u} \in \mathbb{R}^p, \phi_{\text{GAL}_p(\boldsymbol{\Sigma}, \boldsymbol{\mu}, s)}(\mathbf{u}) = \left( \frac{1}{1 + \frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u} - i \boldsymbol{\mu}^T \mathbf{u}} \right)^s.$$

When  $\boldsymbol{\mu} = 0$ , the generalized Laplace distribution is elliptically contoured and is referred to as the *symmetric* generalized Laplace distribution. The elementary properties of the generalized Laplace distribution are discussed by Kozubowski et al. (2013). We list the ones that we consider in the proof of Theorem 1.

**Proposition 3** *If  $\mathbf{z} \sim \text{GAL}_p(\boldsymbol{\Sigma}, \boldsymbol{\mu}, s)$ , we have  $\mathbb{E}(\mathbf{z}) = s\boldsymbol{\mu}$  and  $\text{Cov}(\mathbf{z}) = s(\boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T)$ . Moreover, if  $\boldsymbol{\Sigma}$  is positive definite, the density of  $\mathbf{z}$  is given by*

$$\forall \mathbf{x} \in \mathbb{R}^p, f_{\mathbf{z}}(\mathbf{x}) = \frac{2e^{\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}}}{(2\pi)^{p/2} \Gamma(s) \sqrt{\det \boldsymbol{\Sigma}}} \left( \frac{Q_{\boldsymbol{\Sigma}}(\mathbf{x})}{C(\boldsymbol{\Sigma}, \boldsymbol{\mu})} \right)^{s-p/2} K_{s-p/2}(Q_{\boldsymbol{\Sigma}}(\mathbf{x}) C(\boldsymbol{\Sigma}, \boldsymbol{\mu})), \quad (5)$$

where  $Q_{\Sigma}(\mathbf{x}) = \sqrt{\mathbf{x}^T \Sigma^{-1} \mathbf{x}}$  and  $C(\Sigma, \boldsymbol{\mu}) = \sqrt{2 + \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu}}$ .

**Proposition 4** *Let  $s_1, s_2 > 0, \boldsymbol{\mu} \in \mathbb{R}^p$  and  $\Sigma \in \mathcal{S}_p^+$ . If  $\mathbf{z}_1 \sim \text{GAL}_p(\Sigma, \boldsymbol{\mu}, s_1)$  and  $\mathbf{z}_2 \sim \text{GAL}_p(\Sigma, \boldsymbol{\mu}, s_2)$  are independent random variables, then*

$$\mathbf{z}_1 + \mathbf{z}_2 \sim \text{GAL}_p(\Sigma, \boldsymbol{\mu}, s_1 + s_2). \quad (6)$$

This proposition is a directed consequence of the expression of the characteristic function of the generalized Laplace distribution.

Another appealing property of the multivariate generalized Laplace distribution is that it can be interpreted as an infinite scale mixture of Gaussians with gamma mixing distribution (a property called *Gauss-Laplace transmutation* by Ding and Blitzstein, 2015).

**Proposition 5 (Generalized Gauss-Laplace transmutation)** *Let  $s > 0$  and  $\Sigma \in \mathcal{S}_p^+$ . If  $u \sim \text{Gamma}(s, 1)$  and  $\mathbf{x} \sim \mathcal{N}(0, \Sigma)$  is independent of  $u$ , we have*

$$\sqrt{u}\mathbf{x} \sim \text{GAL}_p(\Sigma, 0, s). \quad (7)$$

For a proof of this result, see Kotz et al. (2001, chap. 6).

To prove Theorem 1, we first study the marginal distribution of the signal term. Following Mattei (2017), we can state the following lemma.

**Lemma 6** *Let  $\mathbf{W}$  be a  $p \times d$  random matrix with i.i.d. columns following a  $\mathcal{N}(0, \phi^{-1}\mathbf{I}_p)$  distribution,  $\mathbf{y} \sim \mathcal{N}(0, \mathbf{I}_d)$  be a Gaussian vector independent from  $\mathbf{W}$ . We obtain*

$$\mathbf{W}\mathbf{y} \sim \text{GAL}_p(2\phi^{-1}\mathbf{I}_p, 0, d/2). \quad (8)$$

**Proof** For each  $k \in \{1, \dots, d\}$  let  $\mathbf{w}_k$  be the  $k$ -th column of  $\mathbf{W}$ ,  $u_k = y_k^2$  and  $\boldsymbol{\xi}_k = y_k \mathbf{w}_k$ . To prove the lemma, we demonstrate that  $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_d$  follow a GAL distribution and use the decomposition

$$\mathbf{W}\mathbf{y} = \sum_{k=1}^d \boldsymbol{\xi}_k.$$

Let  $k \in \{1, \dots, d\}$ . Since  $\mathbf{y}$  is standard Gaussian,  $u_k = y_k^2$  follows a  $\chi^2(1)$  distribution, or equivalently a  $\text{Gamma}(1/2, 1/2)$  distribution. Therefore,  $u_k/2 \sim \text{Gamma}(1/2, 1)$ . Moreover, note that  $\sqrt{u_k} \mathbf{w}_k = |y_k| \mathbf{w}_k = y_k \text{sign}(y_k) \mathbf{w}_k \stackrel{d}{=} y_k \mathbf{w}_k$  since  $|y_k|$  and  $\text{sign}(y_k)$  are independent and  $\text{sign}(y_k) \mathbf{w}_k \stackrel{d}{=} \mathbf{w}_k$ . Therefore, according to Proposition 5, we have

$$\boldsymbol{\xi}_k = \sqrt{\frac{u_k}{2}} \sqrt{2} \mathbf{w}_k \sim \text{GAL}_p(2\phi^{-1}\mathbf{I}_p, 0, 1/2).$$

Since  $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_d$  are i.i.d. and following a  $\text{GAL}_p(2\phi^{-1}\mathbf{I}_p, 0, 1/2)$  distribution, we can use Proposition 4 to conclude that

$$\mathbf{W}\mathbf{y} = \sum_{k=1}^d \boldsymbol{\xi}_k \sim \text{GAL}_p(2\phi^{-1}\mathbf{I}_p, 0, d/2). \quad \blacksquare$$

We now focus on the second term of (1) involving the noise vector.

**Lemma 7** *Let  $\varepsilon_i|\sigma^2 \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_p)$  and  $\sigma^2 \sim \text{Gamma}(a, b)$  then*

$$\varepsilon_i \sim \text{GAL}_p(b^{-1} \mathbf{I}_p, 0, a).$$

**Proof** Again, the Gauss-Laplace transmutation is considered. Indeed, the noise can be written as

$$\varepsilon_i = \sqrt{b\sigma^2} \mathbf{e}_i,$$

where  $\mathbf{e}_i \sim \mathcal{N}(0, b^{-1} \mathbf{I}_p)$ . Therefore, the Gauss-Laplace transmutation allows to conclude. ■

Now that we have proved that both the signal and the noise term follow marginally a generalized Laplace distribution, we use Proposition 4 which ensures that, assuming  $b = \phi/2$ , the sum of the two generalized Laplace random vectors is a generalized Laplace random vector:

$$\mathbf{x}_i \sim \text{GAL}_p(2\phi^{-1} \mathbf{I}_p, 0, a + d/2). \quad (9)$$

Using the expression of the density of the generalized Laplace distribution, we eventually end up with the closed-form expression of the marginal likelihood of Theorem 1.

### 3.3 Choosing hyperparameters

To obtain a closed-form expression of the marginal likelihood, we have shown that it is sufficient to assume that  $b = \phi/2$ . Two hyperparameters remain henceforth to be tuned: the shape parameter of the gamma prior  $a$  and the precision hyperparameter  $\phi$ . We developed data-driven heuristics for this purpose.

A first observation is that, when  $d$  grows,  $\sigma$  is expected to decay because the signal part of the model can be more expressive. This prior information can be distilled into the model by roughly centering the gamma priors on estimates of  $\hat{\sigma}$ . More precisely, our heuristic is to choose  $a$  such that  $\mathbb{E}(\sigma) \propto \hat{\sigma}$  for each  $d$ . In order for  $\phi$  to control the diffusiveness of both the loading matrix and the variance, we specifically made the choice  $a = \hat{\sigma}^2/\phi$ . In our experiments, we chose the ML estimator  $\hat{\sigma} = \sigma_{\text{ML}}$  (which is the mean of the  $p - d$  smallest eigenvalues of the covariance matrix, see Tipping and Bishop, 1999) but more complex estimates may be considered (Passemier et al., 2017).

Regarding the remaining parameter  $\phi$ , we propose a heuristic based on the following statements which can be made regarding the problem of dimension selection:

- overestimation of  $d$  should be preferred to underestimation since losing some information is much more damageable than having a representation not parsimonious enough,
- consequently, the marginal likelihood curve as a function of the dimension should have two distinct phases: a first one when "signal dimensions" are added (before the true value of  $d$ ), and a second one, when "noise dimensions" are added.

Thus, we built a simple heuristic criterion to judge the relevance of a choice of  $\phi$  by the shape of the marginal likelihood curve. First, if the slope of the first part of the curve (before the maximum) is lower than the slope of the second part, this means that this choice leads to underestimation and is therefore discarded. Second, the criterion is equal to the discrete second derivative of the marginal likelihood curve evaluated at the maximum, in



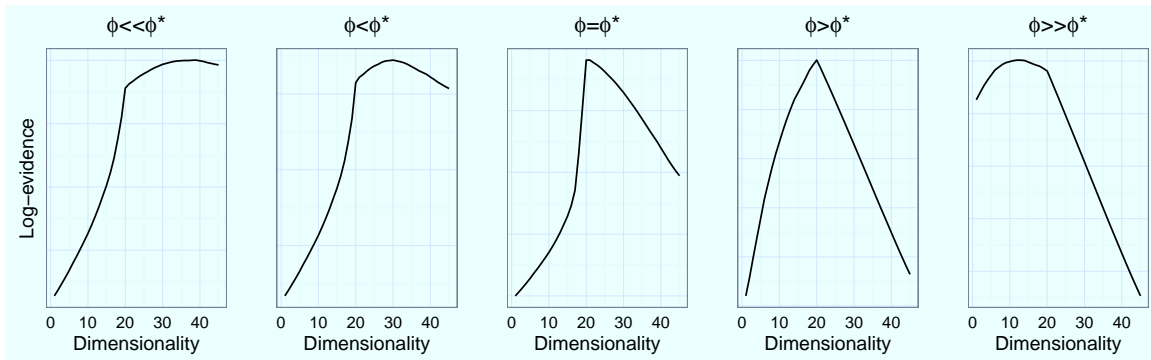


Figure 1: Different shapes of the marginal likelihood curve for growing values of  $\phi$ .  $\phi^*$  corresponds to the maximum of the heuristic criterion that we describe in Subsection 3.3. The true dimensionality is 20.

order to select a hyperparameter leading to a strong distinction between the two phases. This criterion is eventually maximized over a grid of values of  $\phi$ . This scheme for hyperparameter choice is illustrated in Fig. 1 using the simpler simulation scheme described in Subsection 4.2.

## 4. Numerical experiments

In this section, we perform some numerical experiments in order to highlight the main features of the proposed approach and to compare it with state-of-the-art methods.

### 4.1 Simulation scheme

To assess the performance of our algorithm (referred hereafter as ngPPCA or NG, for short), we consider the following simulation scheme in the following experiments. We follow the simulation setup proposed in Bouveyron et al. (2011) based on their isotropic PPCA model. We therefore simulate data sets following the isotropic PPCA model which assumes that the covariance matrix of  $X$  has only two different eigenvalues  $a$  and  $b$  (instead of  $d + 1$  in the PPCA model). In this case, the signal-to-noise ratio (SNR hereafter) is simply defined by

$$\text{SNR} = \frac{ad}{p-d}.$$

In our simulation,  $b$  is set up to 1 and  $a > 1$ , which will control the strength of the signal, varies to explore different signal-to-noise ratios. Then, an orthonormal  $p \times p$  matrix  $\mathbf{Q}$  is drawn uniformly at random. The data is eventually generated according to a centered Gaussian distribution with covariance matrix

$$\mathbf{Q}^T \text{diag}(\overbrace{a, \dots, a}^{d \text{ times}}, \overbrace{1, \dots, 1}^{p-d \text{ times}}) \mathbf{Q}.$$

Finally, the number  $p$  of variables is fixed to 50 in all experiments and the number  $n$  of observations varies in the range  $\{40, 50, 70, 100\}$ .

## 4.2 Introductory examples

We first conduct two small simulations to illustrate the behavior of our algorithm and its difference with the Laplace approximation of Minka (2000). We consider two scenarios: a simple case and a harder and more realistic one.

**Simple scenario** We consider a setup with  $n = 100$  and  $\text{SNR} = 20$ . In this simple scenario, we first illustrate our heuristic for hyperparameter tuning by displaying marginal likelihood curves for different values of  $\phi$  (Fig. 1). The heuristic criterion allows to find the desired shape, leading to a correct dimensionality estimation. A GIF animation displaying all values of the criterion for a large grid of 200 values of  $\phi$  is provided as an online material<sup>1</sup>. On Fig. 2, we compare the results of our algorithm with the Laplace approximation of the marginal likelihood of Minka (2000). In this case, both methods recover the true dimensionality of the data and are very confident with their choice (the posterior probability of the true model is higher than 99% with both approaches). The two curves have a similar shape, in compliance with the expected shape, as detailed in Subsection 3.3.

**Challenging scenario** We now consider a setup with  $n = 40$  and  $\text{SNR} = 20$ . A GIF animation illustrating hyperparameter tuning is provided online<sup>2</sup>. Again, our results are compared with the Laplace approximation (Fig. 3). Regarding our exact approach (left panel), the marginal likelihood curve has an extremely similar shape to the one of the first simulation. This shape is satisfactory, and the maximum of our heuristic criterion actually corresponds to the true dimensionality. Although it also finds the correct dimensionality, the Laplace approximation wrongfully prefers simpler models. More precisely, the top two models chosen by the Laplace approximation are  $\mathcal{M}_{20}$  (with posterior probability 69.3%) and  $\mathcal{M}_{19}$  (with posterior probability 30.7%). In contrast, our algorithm favors  $\mathcal{M}_{20}$  (with posterior probability 86.7%) and  $\mathcal{M}_{21}$  (with posterior probability 13.3%). By preferring overestimation over underestimation, the exact method appears less likely to destroy valuable information, which would be damaging in a dimensionality selection context.

As a summary, those experiments confirm the expected behaviors of NG *vs.* Laplace approximation: in the first scenario ( $n = 100$ ), the asymptotic assumption of the Laplace approximation is much more relevant than in the second setup ( $n = 40$ ). Our method, which does not rely on such an assumption, is much less impacted by the reduction of the sample size.

## 4.3 Benchmark comparison with other dimension selection methods

This section now focuses on the comparison of our methodology with other dimension selection methods. We here consider all possible scenarios with  $n \in \{40, 50, 70, 100\}$  and a SNR grid going from 1.5 to 30 (50 repetitions are made for each case). We compare the performance of our technique based on the normal-gamma prior (NG) with the following four competitors:

- the Laplace approximation of Minka (2000) which is a benchmark Bayesian method for dimension selection,

---

1. <http://pamattei.github.io/animationeasy.gif>

2. <http://pamattei.github.io/animationhard.gif>

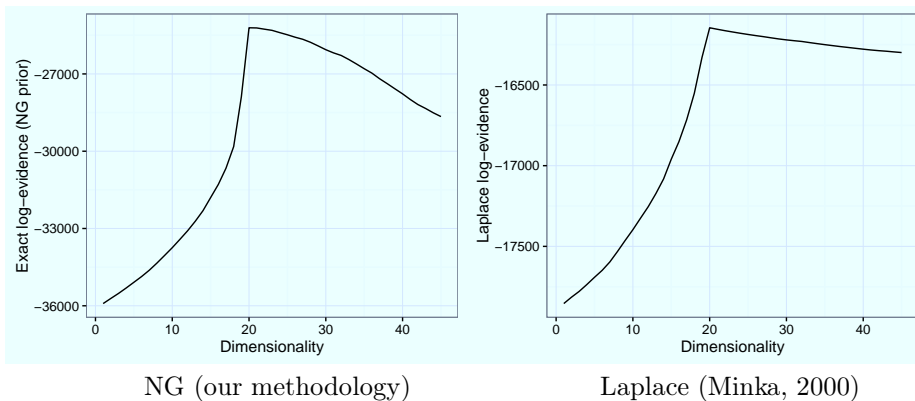


Figure 2: Exact log-evidence for ngPPCA (left) and the Laplace approximation of Minka (2000) (right) for the simpler simulation scenario ( $n = 100$ ). Both curves have the desirable properties detailed in Subsection 3.3 and find the correct dimensionality  $d = 20$ .

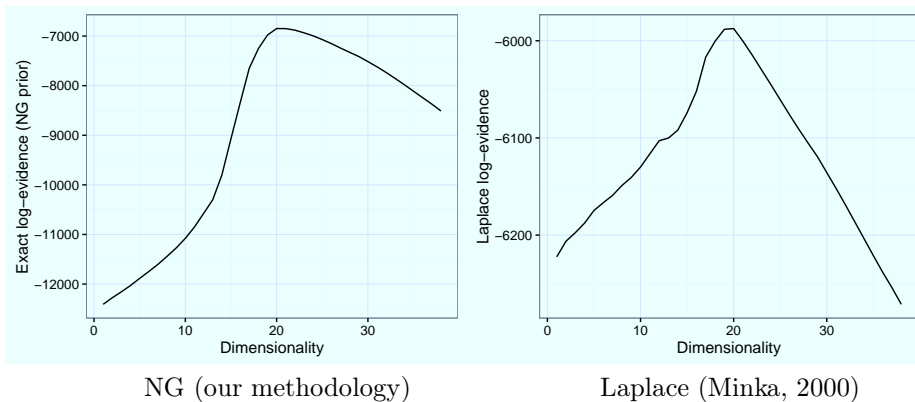


Figure 3: Exact log-evidence for ngPPCA (left) and the Laplace approximation of Minka (2000) (right) for the more challenging simulation scenario ( $n = 40$ ). While both methods select the correct dimensionality  $d = 20$ , the Laplace approximation prefers underestimation which is not a satisfactory behavior. Our exact computation gives a much more acceptable curve.

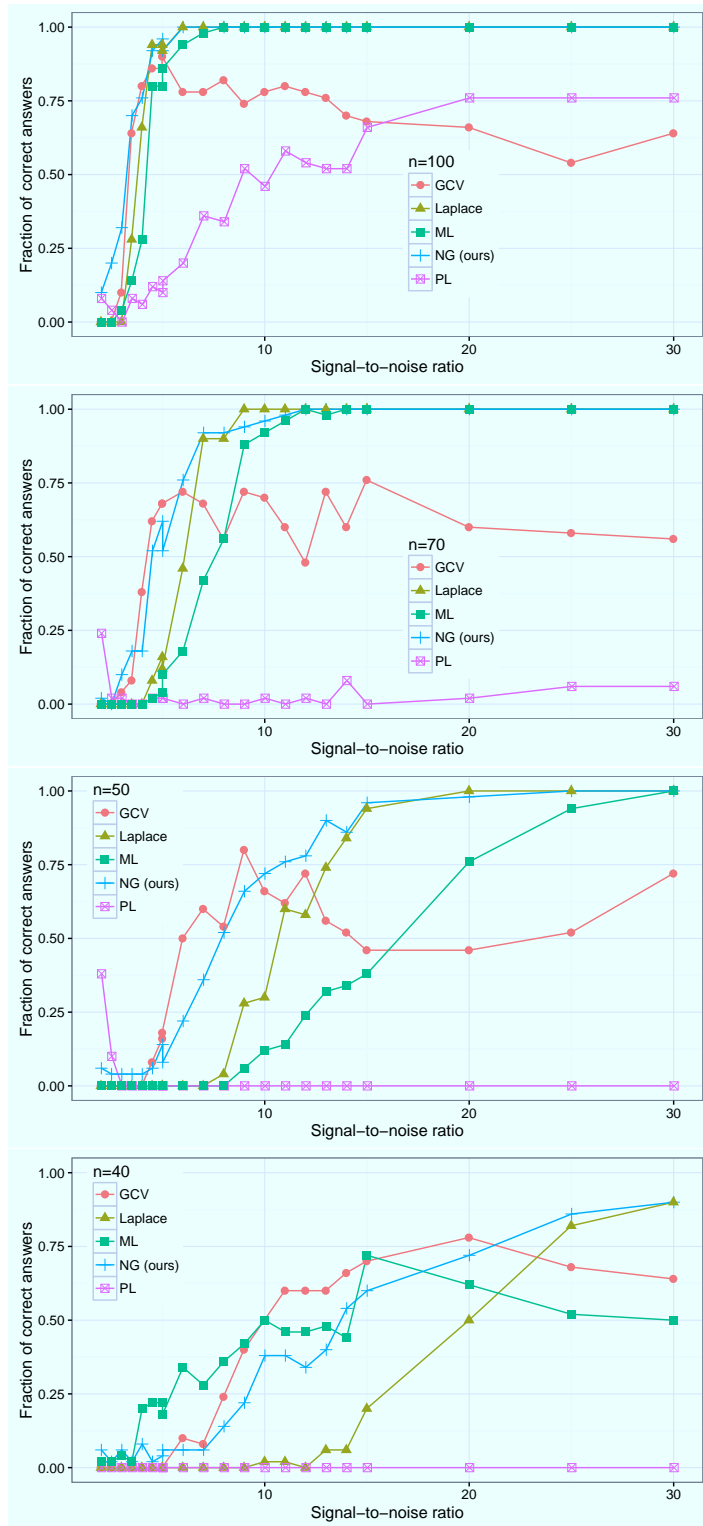


Figure 4: Percentage of correctly estimated dimensions for different sample sizes (averaged 50 replications for each case) of ngPPCA (NG) and its competitors for different signal-to-noise ratios. From top to bottom, the data sample sizes are respectively  $n = 100, 70, 50$  and  $40$ . 11

- the generalized cross-validation approximation (GCV) of Josse and Husson (2012) which is known to give state of the art results in many scenarios (see the vast simulation study of Josse and Husson, 2012)
- the profile likelihood approach (PL) of Zhu and Ghodsi (2006) which represents scree-based techniques and has been very popular in several different contexts (Fogel et al., 2007; Evangelopoulos et al., 2012)
- the ML approach of Bouveyron et al. (2011), which maximizes a non-asymptotic criterion (the likelihood). Notice that is specifically adapted to this simulation scheme and this advantage allow us to consider this technique a gold-standard for the simulated data.

The performance metric that we chose is the percentage of correct answers given by each algorithm, which is a standard measure used in other simulations studies (see e.g. Minka, 2000; Hoyle, 2008; Ulfarsson and Solo, 2008). Results are presented in Fig. 4.

One can first notice that all methods vastly outperform the profile likelihood (PL) approach, which seems not well-suited for small sample sizes. Second, generalized cross-validation gives often satisfactory results, but fails to be competitive with model-based methods (ML, Laplace and NG) when the SNR is high. The ML approach has a good behavior, especially when  $n$  is very small, this is partly explained by the fact that it is designed for this very simulation setup. Finally, our approach (NG), which consistently outperforms the other Bayesian method (the Laplace approximation), is the only method that gives satisfactory results in all settings (high and low SNR, moderate and small  $n$ ).

## 5. Conclusion

PCA is more of a descriptive and exploratory tool than a model. Therefore, no unique dimension selection method should be uniquely preferred – sometimes, very relevant information may actually lie within the *last* PCs (Jolliffe, 2002, section 3.4). However, PCA’s ubiquity in the statistical world makes necessary the search for guidance procedures to help the practitioner choose the number of PCs. This need is even more critical when the data are scarce or particularly expensive. Our work, by deviating from usually adopted asymptotic settings, is a step in that direction. Regarding future work, our exact computation of model posterior probabilities may be used to perform Bayesian model averaging (Hoeting et al., 1999) in predictive contexts. Potential applications could involve principal component regression (Jolliffe, 2002, chap. 8), image denoising (Deledalle et al., 2011), or deep learning (Chan et al., 2015). As a concluding note, this work comes as an illustration that exactly computing the marginal likelihood is sometimes easier than expected. Although both recent asymptotic approximations (Drton and Plummer, 2017) and the MCMC arsenal (Friel and Wyse, 2012) are well-equipped to deal with marginal likelihoods, we argue, like Lin et al. (2009), that finding exact expressions is an important task that should not be deemed untractable too hastily.

## Acknowledgement

Part of this work was conducted while PAM was visiting University College Dublin, funded by the Fondation Sciences Mathématiques de Paris (FSMP).

## References

- H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723, 1974.
- T. Ando. Bayesian factor analysis with fat-tailed factors and its exact marginal likelihood. *Journal of Multivariate Analysis*, 100(8):1717–1726, 2009.
- C. Archambeau and F. Bach. Sparse probabilistic projections. In *Advances in neural information processing systems*, pages 73–80, 2009.
- J. Bai and S. Ng. Determining the number of factors in approximate factor models. *Econometrica*, 70(1):191–221, 2002.
- C. M. Bishop. Bayesian PCA. In *Advances in Neural Information Processing Systems*, pages 382–388, 1999.
- C. Bouveyron, S. Girard, and C. Schmid. High-dimensional data clustering. *Computational Statistics & Data Analysis*, 52(1):502–519, 2007.
- C. Bouveyron, G. Celeux, and S. Girard. Intrinsic dimension estimation by maximum likelihood in isotropic probabilistic PCA. *Pattern Recognition Letters*, 32(14):1706–1713, 2011.
- C. Bouveyron, P. Latouche, and P.-A. Mattei. Bayesian variable selection for globally sparse probabilistic PCA. Technical report, HAL-01310409, 2016.
- R. Bro, K. Kjeldahl, A. K. Smilde, and H. A. L. Kiers. Cross-validation of component models: a critical look at current methods. *Analytical and bioanalytical chemistry*, 390(5):1241–1251, 2008.
- R. B. Cattell. The scree test for the number of factors. *Multivariate behavioral research*, 1(2):245–276, 1966.
- T.-H. Chan, K. Jia, S. Gao, J. Lu, Z. Zeng, and Y. Ma. PCANet: A simple deep learning baseline for image classification? *IEEE Transactions on Image Processing*, 24(12):5017–5032, 2015.
- C.-A. Deledalle, J. Salmon, and A. S. Dalalyan. Image denoising with patch based PCA: local versus global. In *Proceedings of the British Machine Vision Conference*, pages 25.1–25.10, 2011.
- P. Ding and J. K. Blitzstein. Representation for the Gauss-Laplace transmutation. *arXiv preprint arXiv:1510.08765*, 2015.

- M. Drton and M. Plummer. A Bayesian information criterion for singular models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(2):323–380, mar 2017. ISSN 13697412.
- N. Evangelopoulos, X. Zhang, and V. R. Prybutok. Latent semantic analysis: five methodological recommendations. *European Journal of Information Systems*, 21(1):70–86, 2012.
- P. Fogel, S. S. Young, D. M. Hawkins, and N. Ledirac. Inferential, robust non-negative matrix factorization analysis of microarray data. *Bioinformatics*, 23(1):44, 2007.
- N. Friel and J. Wyse. Estimating the evidence—a review. *Statistica Neerlandica*, 66(3):288–308, 2012.
- M. Gavish and D. L. Donoho. The optimal hard threshold for singular values is  $4/\sqrt{3}$ . *IEEE Transactions on Information Theory*, 60(8):5040–5053, 2014.
- E. I. George and R. E. McCulloch. Variable selection via gibbs sampling. *Journal of the American Statistical Association*, 88(423):881–889, 1993.
- A. Hannachi, I. T. Jolliffe, D. B. Stephenson, and N. Trendafilov. In search of simple structures in climate: simplifying EOFs. *International journal of climatology*, 26(1):7–28, 2006.
- J. A. Hoeting, D. Madigan, A. E. Raftery, and C. T. Volinsky. Bayesian model averaging: a tutorial. *Statistical Science*, pages 382–401, 1999.
- P. D. Hoff. Model averaging and dimension selection for the singular value decomposition. *Journal of the American Statistical Association*, 102(478):674–685, 2007.
- H. Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of educational psychology*, 24(6):417, 1933.
- David C Hoyle. Automatic PCA dimension selection for high dimensional data and small sample sizes. *Journal of Machine Learning Research*, 9(Dec):2733–2759, 2008.
- D. A. Jackson. Stopping rules in principal components analysis: a comparison of heuristical and statistical approaches. *Ecology*, 74(8):2204–2214, 1993.
- I. Jolliffe. *Principal component analysis*. Wiley Online Library, 2002.
- I. T. Jolliffe and J. Cadima. Principal component analysis: a review and recent developments. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 374(2065), 2016.
- J. Josse and F. Husson. Selecting the number of components in principal component analysis using cross-validation approximations. *Computational Statistics & Data Analysis*, 56(6):1869–1879, 2012.
- R. E. Kass and A. E. Raftery. Bayes factors. *Journal of the american statistical association*, 90(430):773–795, 1995.

- S. Kotz, T. Kozubowski, and K. Podgórski. *The Laplace distribution and generalizations: a revisit with applications to communications, exconomics, engineering, and finance*. Number 183. Springer Science & Business Media, 2001.
- T. Kozubowski, K. Podgórski, and I. Rychlik. Multivariate generalized laplace distribution and related random fields. *Journal of Multivariate Analysis*, 113:59–72, 2013.
- D. N. Lawley. A modified method of estimation in factor analysis and some large sample results. *Proceedings of the Uppsala Symposium on Psychological Factor Analysis, Uppsala, Sweden*, pages 35–42, 1953.
- S. Lin, B. Sturmfels, and Z. Xu. Marginal likelihood integrals for mixtures of independence models. *Journal of Machine Learning Research*, 10(Jul):1611–1631, 2009.
- P.-A. Mattei. Multiplying a Gaussian matrix by a Gaussian vector (and the Gauss-Laplace transmutation). *arXiv preprint arXiv:1702.02815*, 2017.
- T. P. Minka. Automatic choice of dimensionality for PCA. In *Advances in Neural Information Processing Systems*, volume 13, pages 598–604, 2000.
- K. P. Murphy. Conjugate Bayesian analysis of the Gaussian distribution. *Technical report*, 2007.
- D. Passemier, Z. Li, and J. Yao. On estimation of the noise variance in high dimensional probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(1):51–67, jan 2017. ISSN 13697412.
- K. Pearson. On lines and planes of closest fit to systems of point in space. *Philosophical Magazine*, 2(11):559–572, 1901.
- M. Ringnér. What is principal component analysis? *Nature biotechnology*, 26(3):303–304, 2008.
- G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- C. M. Theobald. An inequality with application to multivariate analysis. *Biometrika*, 62(2):461–466, 1975. ISSN 00063444.
- M. E. Tipping and C. M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622, 1999.
- M. O. Ulfarsson and V. Solo. Dimension estimation in noisy PCA with SURE and random matrix theory. *IEEE Transactions on Signal Processing*, 56(12):5804–5816, 2008.
- S. Wold. Cross-validatory estimation of the number of components in factor and principal components models. *Technometrics*, 20(4):397–405, 1978.
- Y. Zhang and L. El Ghaoui. Large-scale sparse principal component analysis with application to text data. In *Advances in Neural Information Processing Systems*, pages 532–539, 2011.



M. Zhu and A. Ghodsi. Automatic dimensionality selection from the scree plot via the use of profile likelihood. *Computational Statistics & Data Analysis*, 51(2):918–930, 2006.