
Rotation Invariant Householder Parameterization for Bayesian PCA

Rajbir S. Nirwan¹ Nils Bertschinger^{1,2}

Abstract

We consider probabilistic PCA and related factor models from a Bayesian perspective. These models are in general not identifiable as the likelihood has a rotational symmetry. This gives rise to complicated posterior distributions with continuous subspaces of equal density and thus hinders efficiency of inference as well as interpretation of obtained parameters. In particular, posterior averages over factor loadings become meaningless and only model predictions are unambiguous. Here, we propose a parameterization based on Householder transformations, which remove the rotational symmetry of the posterior. Furthermore, by relying on results from random matrix theory, we establish the parameter distribution which leaves the model unchanged compared to the original rotationally symmetric formulation. In particular, we avoid the need to compute the Jacobian determinant of the parameter transformation. This allows us to efficiently implement probabilistic PCA in a rotation invariant fashion in any state of the art toolbox. Here, we implemented our model in the probabilistic programming language Stan and illustrate it on several examples.

1. Introduction

Modern algorithms and computational tools have vastly expanded the scope of Bayesian modeling over the last decades. In particular, Hamiltonian Monte-Carlo (HMC) sampling (see [Betancourt \(2017\)](#) for a thorough introduction) and variational inference ([Bishop, 2006](#)) allow to approximate high-dimensional posterior distributions. In addition, software tools such as probabilistic programming

languages, e.g. Stan ([Carpenter et al., 2017](#)), or libraries for automatic differentiation, e.g. Tensorflow ([Abadi et al., 2015](#)), simplify the implementation. Thanks to these advances, researchers can focus on the statistical modeling as model implementation often requires just a few lines of code. Nevertheless, some models remain challenging. Well known examples include Probabilistic PCA (PPCA) and related factor models which are widely used in exploratory and confirmatory data analysis. These models are non-identifiable, which poses major problems when fitted parameters, e.g. factor loadings, are being interpreted. In this context, the non-identifiability arises from a rotational symmetry of the likelihood model. While identification can be restored by imposing constraints on the factor loadings ([Jöreskog, 1969](#); [Peeters, 2012](#)), maximum likelihood estimation can be biased by the imposed constraints ([Millsap, 2001](#)).

Here, we consider PPCA from a Bayesian perspective. In this context, the rotational symmetry gives rise to a complicated posterior with continuous subspaces of equal density. Continuous symmetries can severely reduce sampling efficiency as many samples are required in order to explore the corresponding subspaces. Furthermore, interpretation of the sampled parameters becomes impossible as marginal posterior averages are meaningless when many different parameter combinations give rise to the same likelihood. This is akin to the discrete label switching symmetry observed in mixture models. In this case, the parameters of the individual mixture components cannot be averaged as the arbitrary component labels between different samples could be switched. Thus, a major advantage of Bayesian modeling, namely the ability to access estimation uncertainty, cannot be utilized as posterior means and variances not just reflect variation within, but also between components.

Therefore, it is desirable to remove the rotational symmetry of PPCA. Formally, the symmetry arises due to the choice of coordinate system for the principal components being arbitrary. In particular, a rotation of the coordinate axes leaves the model likelihood unchanged. Mathematically, the rotation symmetry is made explicit by formulating the model in terms of an orthogonal matrix (transforming into a fixed coordinate system) and a diagonal matrix containing the explained variances. To this end, we employ the singular value decomposition (SVD) and parameterize the Stiefel manifold of orthogonal matrices in terms of non-

¹Department of Computer Science, Goethe University, Frankfurt, Germany ²Frankfurt Institute for Advanced Studies, Frankfurt, Germany. Correspondence to: Rajbir S. Nirwan <nirwan@fias.uni-frankfurt.de>, Nils Bertschinger <bertschinger@fias.uni-frankfurt.de>.

redundant unconstrained parameters. The mapping between parameters and orthogonal matrices is provided in terms of a sequence of Householder transformations.

1.1. Related work

Special sampling methods for inference on statistical manifolds have been developed, in which the geometric structure of the manifold is respected and geodesic trajectories are traced out with respect to the metric of the manifold (Byrne & Girolami, 2013). In general, the geometry will be non-Euclidean, e.g. spherical for unit vectors. Also the geometry of the Stiefel manifold of orthogonal matrices can be handled in this fashion. Unfortunately, as the geometry needs to be build into the sampling algorithm, this approach requires substantive work. This also applies to approaches using analytic results for the matrix von Mises-Fisher distribution on the Stiefel manifold. Exploiting conditional conjugacy (Hoff, 2009) and (Šmídl & Quinn, 2007) have employed this distribution to implement Gibbs sampling and variational inference for Bayesian PPCA respectively. An alternative, which allows to utilize general purpose tool boxes, is to reparameterize the manifold in terms of unconstrained parameters¹. Pourzanjani et al. (2017) achieved this via Givens rotations. They also demonstrate that the model can then be implemented in Stan without any changes to the underlying sampling routine and unrestricted by conjugacy requirements.

Ideally, reparameterizing a Bayesian model should not change the joint density defined by the model. In case of the transformation via Givens rotations, the parameter density corresponding to a uniform measure on the Stiefel manifold is not known. Thus, by change of measure the density needs to be corrected by the Jacobian determinant, which poses a major computational bottleneck in high-dimensional models. Here, we propose a different parameterization in terms of Householder transformations. In this case, results from random matrix theory allow us to obtain the induced parameter density corresponding to a Gaussian prior on the original parameters of PPCA, which include the rotational symmetry. In section 2 we recall the necessary mathematical results to do so. Then, we illustrate the resulting model on several examples (section 3 and 4) before we conclude in section 5.

2. Background

Here, we quickly review PPCA and explain why the model is rotationally symmetric, and how the symmetry can be

¹In general, such a reparameterization cannot be achieved globally. Instead, the non-Euclidean geometry gives rise to singularities in the parameter transformation, e.g. at the north pole when mapping a sphere onto a plane, or a different mapping, i.e. chart, needs to be employed at different points.

removed via singular value decomposition (SVD). Finally, we discuss some results from random matrix theory in order to construct a suitable prior on the transformed parameters.

2.1. Probabilistic Principal Component Analysis (PPCA)

PPCA (Tipping & Bishop, 1999; Bishop, 2006) relates a D -dimensional observation \mathbf{y} to a Q -dimensional latent vector \mathbf{x} via a linear mapping $\mathbf{W} \in \mathbb{R}^{D \times Q}$

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad (2.1)$$

where $\boldsymbol{\mu}$ allows the model to have a non-zero mean and $\boldsymbol{\epsilon}$ states the variance not explained by the first two terms as noise. PPCA assumes a standard Gaussian distribution on the latent space, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and a zero mean Gaussian noise with variance σ^2 for all dimensions, i.e. $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. The likelihood for \mathbf{y} given \mathbf{W} , $\boldsymbol{\mu}$ and σ^2 takes the form

$$p(\mathbf{y}|\mathbf{W}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}). \quad (2.2)$$

For N observations, denoted as $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N)^T \in \mathbb{R}^{N \times D}$, equation (2.1) becomes

$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \boldsymbol{\epsilon}, \quad (2.3)$$

where $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times Q}$ and $\boldsymbol{\epsilon} \in \mathbb{R}^{N \times D}$. The likelihood of \mathbf{Y} given \mathbf{W} then becomes

$$\begin{aligned} p(\mathbf{Y}|\mathbf{W}) &= \prod_{n=1}^N \mathcal{N}(\mathbf{Y}_{n,:} | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}), \\ \ln p(\mathbf{Y}|\mathbf{W}) &= \frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1} \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^T), \end{aligned} \quad (2.4)$$

where $\mathbf{K} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$ and $\tilde{\mathbf{Y}}$ are the centered observations, i.e. $\tilde{\mathbf{Y}}_{n,:} = \mathbf{Y}_{n,:} - \boldsymbol{\mu}$. This expression can be maximized analytically and yields the solution (Tipping & Bishop, 1999)

$$\begin{aligned} \boldsymbol{\mu}_{\text{ML}} &= \frac{1}{N} \sum_n \mathbf{Y}_{n,:}, \\ \mathbf{W}_{\text{ML}} &= \mathbf{U} (\boldsymbol{\Lambda} - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}, \end{aligned} \quad (2.5)$$

where $\mathbf{U} \in \mathbb{R}^{D \times Q}$ is an orthogonal matrix containing the principal eigenvectors of the empirical covariance matrix $\frac{1}{N} \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^T$, $\boldsymbol{\Lambda} \in \mathbb{R}^{Q \times Q}$ is a diagonal matrix containing the corresponding eigenvalues $(\lambda_1, \dots, \lambda_Q)$ of $\frac{1}{N} \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^T$ on the diagonal and $\mathbf{R} \in \mathbb{R}^{Q \times Q}$ is an arbitrary rotation matrix. The maximum likelihood estimation for σ^2 is given by

$$\sigma_{\text{ML}}^2 = \frac{1}{D - Q} \sum_{i=Q+1}^D \lambda_i, \quad (2.6)$$

which is basically the variance of \mathbf{Y} not picked up by the model. By projecting the data points \mathbf{y} into the latent space, a representation of reduced dimensionality $Q < D$ can be obtained. The posterior mean of the latent vectors is given by (Bishop, 2006)

$$\mathbb{E}[\mathbf{x}|\mathbf{y}] = \left(\tilde{\mathbf{W}}_{\text{ML}}^T \mathbf{W}_{\text{ML}} + \sigma^2 \mathbf{I} \right)^{-1} \tilde{\mathbf{W}}_{\text{ML}}^T (\mathbf{y} - \boldsymbol{\mu}_{\text{ML}}). \quad (2.7)$$

Note that the rotation symmetry in the likelihood caused by \mathbf{R} in (2.5) makes the likelihood non-unique. Two different solutions \mathbf{W} and $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{R}$ will lead to the same likelihood, since $\tilde{\mathbf{W}}\tilde{\mathbf{W}}^T = \mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T$. For this reason, numerical optimization algorithms often converge to different results when started from different initial conditions. While model predictions are unaffected by this non-uniqueness, interpretation of parameters and the projected latent vectors becomes impossible. Especially the latter is problematic if PPCA is employed as a pre-processing step to reduce the dimensionality of the data and further analysis might be sensitive to the arbitrary choice of latent space rotation. Figure 2.1 (left side) shows the maximum likelihood result for 1000 different initial conditions for \mathbf{W} . One can clearly see the rotation symmetry.

PPCA is closely related to classical PCA. In particular, as shown by Tipping & Bishop (1999), the classical solution is recovered by letting $\sigma^2 \rightarrow 0$. In this case, the dimensionality reduction is achieved by the linear projection

$$\mathbf{x} = \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T (\mathbf{y} - \boldsymbol{\mu}_{\text{ML}}) \quad (2.8)$$

where $\boldsymbol{\mu}_{\text{ML}}$, $\boldsymbol{\Lambda}$ and \mathbf{U} are defined as in (2.5). The eigenvalues λ_q are interpreted as the variance explained by the q -th principal component $\mathbf{U}_{:,q}$ in this context. From (2.7) the classical solution is obtained by assuming vanishing noise variance σ^2 and fixing the latent space rotation of the maximum likelihood solution at $\mathbf{R} = \mathbf{I}$.

Factor analysis is closely related to PPCA as well. In this case, PPCA is slightly generalized by assuming that the noise is distributed as $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$ with diagonal covariance matrix $\boldsymbol{\Psi}$. In particular, the likelihood has the same rotational symmetry as PPCA. Here, we just note that all our derivations in the following sections apply equally to this model.

2.2. Bayesian Approach to PPCA

The fully Bayesian way to solve (2.4) would be to impose prior distributions $p(\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$ and solve for the posterior

$$p(\boldsymbol{\mu}, \mathbf{W}, \sigma^2 | \mathbf{Y}) = \frac{p(\mathbf{Y} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) p(\boldsymbol{\mu}, \mathbf{W}, \sigma^2)}{p(\mathbf{Y})}, \quad (2.9)$$

which is not tractable anymore. Therefore, we have to resort to other techniques, e.g. sampling the posterior. For

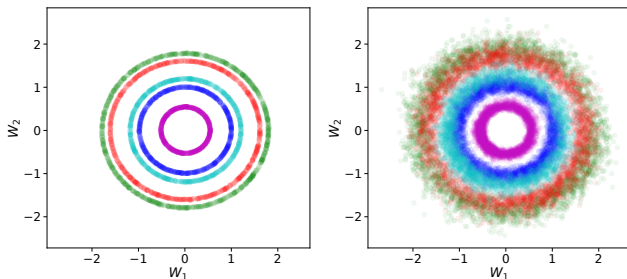


Figure 2.1. Data are created by sampling 50 5-dimensional points from the standard normal distribution and map them via a random linear mapping $\mathbf{W} \in \mathbb{R}^{5 \times 2}$ to the data-space \mathbf{Y} . Elements of \mathbf{W} are drawn from a standard normal distribution. Left side: Maximum likelihood solution for 1000 runs with different initial conditions of \mathbf{W} . Right side: 4000 samples from the posterior distribution. Different colors show different rows of \mathbf{W} .

a rotation invariant prior $p(\mathbf{W})$ (e.g. an isotropic Gaussian) the posterior will have the same rotation symmetry as well. Thus, the posterior has a complicated shape with a continuous subspace of equal density posing a considerable challenge for sampling algorithms.

Figure 2.1 (right side) shows 4000 posterior samples. The likelihood is given by (2.4) and the prior on \mathbf{W} is a standard normal. We can clearly see the rotation symmetry again. Samples were drawn using the NUTS sampler of Stan (Carpenter et al., 2017), which is able to fully explore the challenging posterior. Yet, in higher dimensions $D > Q \gg 1$ the problem becomes more severe and a lot of samples are required to fully explore the symmetric solutions. Furthermore, statistics of the posterior samples, e.g. computing the mean or variance, do not have any significance in this case.

The question arises: Can we somehow take out the rotation symmetry without changing the model (i.e. likelihood or the mass of the posterior distribution)? The answer is yes: We can achieve this by reparameterizing \mathbf{W} . That way the results are identified uniquely. The aim is not to change the model by the reparameterization. Therefore, we need specific corresponding distributions for the parameters on the reparameterized model. In the next sections we explain how to do that.

2.3. Singular Value Decomposition

It is well known, that every positive-definite matrix \mathbf{S} can be diagonalized by an orthogonal transformation \mathbf{U} as $\mathbf{S} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$, where $\boldsymbol{\Lambda}$ contains the eigenvalues of \mathbf{S} on the diagonal and the columns of \mathbf{U} are the corresponding normalized eigenvectors.

This is not true anymore for other matrices, but there exists a more general decomposition called the singular value

decomposition (SVD). A matrix $\mathbf{A} \in \mathbb{R}^{D \times Q}$ can be decomposed as

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T, \quad (2.10)$$

where $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_D)^T$, $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_Q)$ are orthogonal matrices and $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_Q) \in \mathbb{R}^{D \times Q}$ is a diagonal matrix containing the singular values. When \mathbf{A} has rank $P < Q$, only P of them will be non-zero. It is easy to show that the square of the diagonal elements $(\sigma_1^2, \dots, \sigma_P^2)$ are the non-zero eigenvalues of $\mathbf{A}\mathbf{A}^T$ and $\mathbf{A}^T\mathbf{A}$, which are similar matrices, and all \mathbf{u}_j 's and \mathbf{v}_j 's obey the relation

$$\begin{aligned} \mathbf{A}\mathbf{v}_j &= \sigma_j\mathbf{u}_j, & \mathbf{A}^T\mathbf{u}_j &= \sigma_j\mathbf{v}_j, & \forall j &= 1, \dots, P \\ \mathbf{A}\mathbf{v}_j &= \mathbf{0}, & \mathbf{A}^T\mathbf{u}_j &= \mathbf{0}, & \forall j &> P \end{aligned} \quad (2.11)$$

We will decompose our mapping $\mathbf{W} \in \mathbb{R}^{D \times Q}$ via SVD as $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. The likelihood only depends on the outer product of \mathbf{W} . Since the outer product does not depend on \mathbf{V} ,

$$\mathbf{W}\mathbf{W}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T = \mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^T, \quad (2.12)$$

we can neglect \mathbf{V}^2 , i.e., we assume it to be the identity matrix. That fixes the coordinate frame of the latent space and takes out the rotation symmetry.

By reparameterizing \mathbf{W} by \mathbf{U} and $\mathbf{\Sigma}$, we obtain a unique representation³. However, in the fully Bayesian case, we have a prior on \mathbf{W} . After reparameterization by \mathbf{U} and $\mathbf{\Sigma}$, we have to make a Jacobian adjustment to the probability density, since otherwise the prior density mass on \mathbf{W} will change as well. In the next sections, we consider a zero mean Gaussian prior on \mathbf{W} , and how to adjust the prior on \mathbf{U} and $\mathbf{\Sigma}$, such that we do not change the distribution on \mathbf{W} .

Similarly, SVD provides a numerically stable way for classical PCA. In this case, the data matrix \mathbf{Y} is decomposed as $\mathbf{Y} = \mathbf{U}_{PCA}\mathbf{\Sigma}_{PCA}\mathbf{V}_{PCA}^T$ and the desired dimensionality reduction is achieved by the linear map $\mathbf{\Sigma}_{PCA}\mathbf{U}_{PCA}^T$. Again, we see that the matrix of right singular vectors \mathbf{V}_{PCA} , corresponding to an arbitrary choice of latent coordinate rotation, is dropped.

2.4. Random Matrix Theory, Haar Measure

SVD decomposes our \mathbf{W} into an orthogonal matrix \mathbf{U} , a diagonal matrix $\mathbf{\Sigma}$ and another orthogonal matrix \mathbf{V} , that we are neglecting in our analysis, since the outer product $\mathbf{W}\mathbf{W}^T$ does not depend on \mathbf{V} .

²Note, \mathbf{V} is precisely the rotation symmetry that we want to get rid of.

³In general, the SVD is not unique. While the singular values can be ordered to ensure uniqueness of $\mathbf{\Sigma}$, the left singular vectors \mathbf{U} are only determined up to the sign. As explained later, we enforce a sign convention on them.

For a standard Gaussian prior on \mathbf{W} , i.e. all elements of \mathbf{W} are zero mean Gaussian with unit variance, $\mathbf{W}\mathbf{W}^T$ has a Wishart distribution and the following theorem holds (James & Lee, 2014):

Theorem 1 *Let the entries of $\mathbf{W} \in \mathbb{R}^{D \times Q}$ be i.i.d. Gaussian with zero mean and unit variance. The joint probability density of the ordered strictly positive eigenvalues of the Wishart matrix $\mathbf{W}^T\mathbf{W}$, $\lambda_1 \geq \dots \geq \lambda_Q$, equals*

$$p(\boldsymbol{\lambda}) = ce^{-\frac{1}{2}\sum_{q=1}^Q \lambda_q} \prod_{q=1}^Q \left(\lambda_q^{\frac{D-Q-1}{2}} \prod_{q'=q+1}^Q |\lambda_q - \lambda_{q'}| \right), \quad (2.13)$$

where $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_Q)$ and c is a constant that depends on D and Q .

Note that the non-zero eigenvalues of $\mathbf{W}^T\mathbf{W}$ are the same as for $\mathbf{W}\mathbf{W}^T$, and therefore they have the same probability density function. The singular values σ_i of \mathbf{W} are the square roots of λ_i 's, $\sigma_i = \sqrt{\lambda_i}$. Thus, by change of measure we get⁴

$$\begin{aligned} p(\sigma_1, \dots, \sigma_Q) &= \\ ce^{-\frac{1}{2}\sum_{q=1}^Q \sigma_q^2} \prod_{q=1}^Q \left(\sigma_q^{D-Q-1} \prod_{q'=q+1}^Q |\sigma_q^2 - \sigma_{q'}^2| \right) \prod_{q=1}^Q 2\sigma_q, \end{aligned} \quad (2.14)$$

where the last product in the above term is the Jacobian correction. Provided this density function as the prior distribution on the singular value matrix $\mathbf{\Sigma}$, we can easily sample from its posterior.

To calculate the prior distribution for \mathbf{U} , we need to dig deeper into random matrix theory. As mentioned earlier, a standard Gaussian prior on \mathbf{W} gives rise to a Wishart distribution on $\mathbf{W}\mathbf{W}^T$. Equation (2.12) shows, that a Wishart matrix can be decomposed into a product of an orthogonal matrix \mathbf{U} and a diagonal matrix $\mathbf{\Sigma}$. Furthermore, it is known that the eigenvectors \mathbf{U} are distributed uniformly in the space of orthogonal matrices⁵. The set of orthogonal matrices is known as the Stiefel manifold $\mathcal{V}_{Q,D}$

$$\mathcal{V}_{Q,D} = \left\{ \mathbf{U} \in \mathbb{R}^{D \times Q} \mid \mathbf{U}^T\mathbf{U} = \mathbf{I} \right\}. \quad (2.15)$$

The dimension of this manifold is $DQ - \frac{1}{2}Q(Q+1)$, accounting for the fact that the orthogonality constraint reduces the number of independent degrees of freedom. The

⁴Given a probability density $p_x(x)$ on x and an invertible map g , so that $y = g(x)$, the density $p_y(y)$ on y is given by $p_y(y) = p_x(g^{-1}(y)) \left| \det \left(\frac{dg^{-1}(y)}{dy} \right) \right|$. Here the absolute determinant of the Jacobian accounts for the change of volume under g .

⁵Bai et al. (2007) and Uhlig (1994) mention that the eigenvector matrix of a Wishart distribution is Haar-distributed, which is also true for singular Wishart matrices, but we could not find a proof anywhere.

Stiefel manifold can be equipped with a uniform measure. This measure is an example of a *Haar* measure, as it is invariant under the action of the orthogonal group $O(D) \simeq \mathcal{V}_{D,D}$, i.e.

$$p(\mathbf{U}) = p(\mathbf{R}\mathbf{U}) \quad \forall \mathbf{R} \in O(D). \quad (2.16)$$

To proceed, we need to find an unconstrained parameterization for orthogonal matrices along with a density on the parameters, such that the resulting matrix is Haar distributed. Shepard et al. (2014); Pourzanjani et al. (2017) and Mezzadri (2007) suggest ways to do that. The procedure is explained in detail in the next section.

2.5. QR-decomposition and Householder Transformations

Given a matrix $\mathbf{Z} \in \mathbb{R}^{D \times Q}$ the (thin) QR-decomposition decomposes \mathbf{Z} into an orthogonal matrix $\mathbf{Q} \in \mathcal{V}_{Q,D}$ and an upper triangular matrix $\mathbf{R} \in \mathbb{R}^{Q \times Q}$. If the elements of \mathbf{Z} are i.i.d Gaussian with mean zero and unit variance, \mathbf{Q} is Haar distributed (Mezzadri, 2007). Note that this is only the case if a unique QR-decomposition is used. In practice, this is usually achieved by enforcing the convention that the diagonal elements of \mathbf{R} are positive.

To compute the QR-decomposition of \mathbf{Z} , the so-called Householder Transformations \mathbf{H} can be used. These transformations are reflections on the plane spanned by a vector $\mathbf{v}_n \in \mathbb{R}^n$. To decompose a D -by- Q matrix, Q of such transformations are needed and the resulting orthogonal matrix \mathbf{Q} can be written as

$$\mathbf{Q} = \mathbf{H}_D(\mathbf{v}_D)\mathbf{H}_{D-1}(\mathbf{v}_{D-1})\dots\mathbf{H}_{D-Q+1}(\mathbf{v}_{D-Q+1}). \quad (2.17)$$

To construct \mathbf{H}_n , we define $\tilde{\mathbf{H}}_n(\mathbf{v}_n)$ as

$$\tilde{\mathbf{H}}_n(\mathbf{v}_n) = -\text{sgn}(\mathbf{v}_{n1}) (\mathbf{I} - 2\mathbf{u}_n\mathbf{u}_n^T) \in \mathbb{R}^{n \times n}, \quad (2.18)$$

where

$$\mathbf{u}_n = \frac{\mathbf{v}_n + \text{sgn}(\mathbf{v}_{n1})\|\mathbf{v}_n\|\mathbf{e}_1}{\|\mathbf{v}_n + \text{sgn}(\mathbf{v}_{n1})\|\mathbf{v}_n\|\mathbf{e}_1\|} \quad (2.19)$$

and construct \mathbf{H}_n by

$$\mathbf{H}_n = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{H}}_n \end{pmatrix}. \quad (2.20)$$

The algorithm for the QR-decomposition constructs a suitable sequence of vectors $\mathbf{v}_D, \dots, \mathbf{v}_{D-Q+1}$ based on the entries of the successively rotated columns of \mathbf{Z} . The choice of the sign fulfills the requirements mentioned earlier, i.e. the diagonal elements of \mathbf{R} are positive.

Here, we are only interested in the resulting orthogonal matrix \mathbf{Q} . In particular, by randomly drawing vectors \mathbf{v} , we obtain a random orthogonal matrix via (2.17). The following theorem tells us which distribution the \mathbf{v} s need to have in order to get an orthogonal matrix distributed according to the Haar measure (Mezzadri, 2007).

Theorem 2 *Let $\mathbf{v}_D, \mathbf{v}_{D-1}, \dots, \mathbf{v}_1$ be uniformly distributed on the unit spheres $\mathbb{S}^{D-1}, \dots, \mathbb{S}^0$ respectively, where \mathbb{S}^{n-1} is the unit sphere in \mathbb{R}^n . Furthermore, let $\mathbf{H}_n(\mathbf{v}_n)$ be the n -th Householder transformation as defined in equation (2.20) The product*

$$\mathbf{Q} = \mathbf{H}_D(\mathbf{v}_D)\mathbf{H}_{D-1}(\mathbf{v}_{D-1})\dots\mathbf{H}_1(\mathbf{v}_1) \quad (2.21)$$

is a random orthogonal matrix with distribution given by the Haar measure on $O(D)$.

A corresponding draw from the Stiefel manifold $\mathcal{V}_{Q,D}$ can be obtained by just taking the first Q columns of \mathbf{Q} . Alternatively, this is achieved by drawing vectors $\mathbf{v}_D, \dots, \mathbf{v}_{D-Q+1}$ uniformly from the respective unit sphere and construct the unitary matrix from the corresponding Householder transformations. By the theorem above, this gives the Haar measure on $\mathcal{V}_{Q,D}$ as the Householder transformations $\mathbf{H}_{D-Q}, \dots, \mathbf{H}_1$ only effect the columns $D-Q, \dots, 1$ by (2.20).

3. Unique PPCA

Here, we connect the results from the previous section to obtain a model for PPCA without the latent space symmetry that plaques the original formulation. We start by sampling uniformly from the unit spheres \mathbb{S}^{n-1} . This is most easily achieved by drawing an i.i.d. standard Gaussian vector of dimension n and normalizing its length. Note that this spends an additional parameter compared to the dimensionality $n-1$ of the unit sphere \mathbb{S}^{n-1} . Yet, the n -dimensional standard Gaussian distribution is easy to sample from and the vector length is sufficiently constrained by the Gaussian prior such that the sampler can move around the unit sphere effectively. The same parameterization is employed, for example, by Stan in order to support a `unit_vector` data type (Stan Development Team, 2018). Furthermore, the Householder transformation does not require the vectors \mathbf{v} to be of unit length as they are normalized anyways via (2.19). Thus, we construct the orthogonal matrix \mathbf{U} distributed with the Haar measure by successively transforming standard Gaussian random vectors.

Next, we sample the ordered singular values from the joint distribution (2.14). Again, this can be accomplished by transforming the ordered vector $(\sigma_1, \dots, \sigma_Q)$ to an unconstrained space and correcting the joint density by the Jacobian determinant of the transformation. We refer to Stan Development Team (2018) for details of this transformation. Finally, from the orthogonal matrix \mathbf{U} and the diagonal matrix $\mathbf{\Sigma}$, we construct \mathbf{W} by $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}$, which we then use in the likelihood of the PPCA model. Overall, we obtain the

following generative model:

$$\begin{aligned}
 \mathbf{v}_D, \dots, \mathbf{v}_{D-Q+1} &\sim \mathcal{N}(0, \mathbf{I}) \\
 \boldsymbol{\sigma} &\sim p(\boldsymbol{\sigma}) \propto \text{eq. (2.14)} \\
 \boldsymbol{\mu} &\sim p(\boldsymbol{\mu}), \text{ e.g. a broad Gaussian} \\
 \mathbf{U} &= \prod_{q=1}^Q \mathbf{H}_{D-q+1}(\mathbf{v}_{D-q+1}) \\
 \boldsymbol{\Sigma} &= \text{diag}(\boldsymbol{\sigma}) \\
 \mathbf{W} &= \mathbf{U}\boldsymbol{\Sigma} \\
 \sigma_{\text{noise}} &\sim p(\sigma_{\text{noise}}) \\
 \mathbf{Y} &\sim \prod_{n=1}^N \mathcal{N}(\mathbf{Y}_{n,:} | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma_{\text{noise}}^2 \mathbf{I}).
 \end{aligned}$$

Note that this model defines the same distribution as the corresponding model with a standard Gaussian prior $\mathbf{W} \sim \mathcal{N}(0, \mathbf{I})$. In both cases, the sampling distribution is governed by the Wishart distributed matrix $\mathbf{W}\mathbf{W}^T$, even though the distribution on \mathbf{W} is actually different. We implemented both models in the probabilistic programming language Stan. The code for the simulations is available on Github: <https://github.com/RSNirwan/HouseholderBPCA>.

Compared to previous approaches based on parameterizing the Stiefel manifold in terms of Givens rotations (Pourzajani et al., 2017), our model has the following advantages: First, our Householder parameters \mathbf{v} are unconstrained, in contrast to the angular parameters of Givens rotations where the sampler might hit the boundary of the space. Secondly, we avoid the computationally demanding computation of the Jacobian determinant (Shepard et al., 2014). Compared to other approaches employing the matrix von Mises-Fisher distribution (Hoff, 2009; Šmídl & Quinn, 2007), our model has the following advantages: First, we do not require conditional conjugacy allowing non-linear extensions of Bayesian PPCA as illustrated in section 4. Secondly, we do not need to resort to rejection sampling or variational approximations. Similar to these approaches, our parameterization introduces a combinatorial symmetry by the sign ambiguity of the SVD. This is akin to the label switching problem in Gaussian mixture models which usually poses few problems with the sampler simply getting stuck in a specific mode. Here, in order to compare results across different modes, we postprocess each sample such that the first entry of each column of \mathbf{U} is made positive. In the following, we illustrate our model on some data sets and discuss possible extensions to non-linear models where similar symmetries arise.

3.1. Model Comparison

3.1.1. SYNTHETIC DATASET

Here, we build our own synthetic dataset with known parameters and the goal is to reconstruct the parameter values.

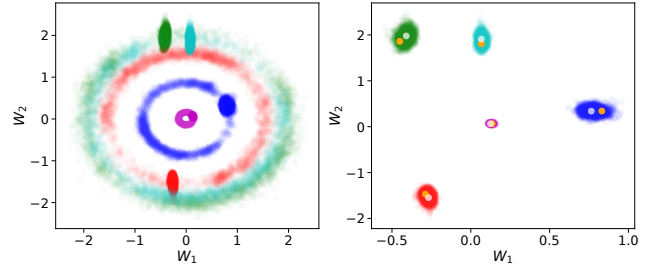


Figure 3.1. Results for synthetic dataset. Left: In the background, we see samples from the posterior of the standard parameterization. Samples from the \mathbf{U} and $\boldsymbol{\Sigma}$ mapped to $\mathbf{W} = \mathbf{U}\boldsymbol{\Sigma}$ are shown in a darker color. Right: Comparison of the suggested model to classical PCA. White dots are the classical PCA solution and orange dots are the true values. 4000 samples for each row of \mathbf{W} are shown in different colors.

For $(N, D, Q) = (150, 5, 2)$ we sample $\mathbf{X} \in \mathbb{R}^{N \times D}$ from a standard normal distribution and construct \mathbf{W} by $\mathbf{W} = \mathbf{U}\boldsymbol{\Sigma} \in \mathbb{R}^{D \times Q}$, where \mathbf{U} is sampled from the Stiefel manifold with Haar measure and we specify $\boldsymbol{\Sigma} = \text{diag}(\sigma_1, \sigma_2)$, where $(\sigma_1, \sigma_2) = (3.0, 1.0)$. Then, we get the observation $\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon}$ denotes the noise sampled from a zero mean Gaussian with a standard deviation of 0.01.

The left plot in Figure 3.1 compares the Bayesian inference for the standard model, where we directly assume a standard Gaussian prior on \mathbf{W} with our suggested model, where we parameterize \mathbf{W} by \mathbf{U} and $\boldsymbol{\Sigma}$ and infer the posterior distribution of both. As expected, the standard model has the rotational symmetry in \mathbf{W} , whereas our model takes out the rotation⁶. On the right side (Figure 3.1), we compare the posterior distribution of our model with the classical PCA solution and as expected with the low observation noise, the solutions are quite similar. In contrast to classical PCA, the Bayesian approach provides a distribution for our parameters, accessing estimation uncertainty of the solution as well.

Figure 3.2 shows the posterior distribution of $\boldsymbol{\Sigma}$, that we fixed to $\text{diag}(3.0, 1.0)$. As we can see, the true values are recovered quite well.

3.1.2. BREAST CANCER WISCONSIN DATASET

We tested the model on the Breast Cancer Wisconsin dataset as well. The dataset was downloaded from the Python toolbox scikit-learn (Pedregosa et al., 2011) and contains 569 labeled datapoints with 30 features. We neglect the labels and take the standardized 569×30 matrix as the input to our model. For visualization purposes, we again map the data to $Q = 2$. Figure 3.3 shows the performance of different models. In the left plot, we see the posterior samples for the

⁶The additional sign ambiguity of the columns of \mathbf{U} is removed afterwards by post processing each sample.

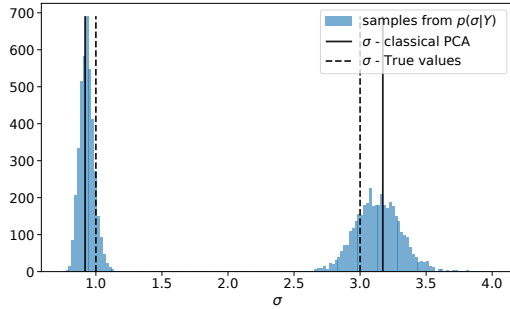


Figure 3.2. Histogram of the samples from the posterior distribution of $\Sigma = \text{diag}(\sigma_1, \sigma_2)$. The vertical solid black lines are the classic PCA solution (3.17, 0.92) and dashed lines are the true values (3.0, 1.0), which generated the data.

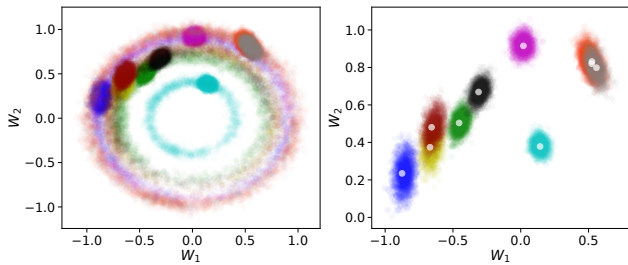


Figure 3.3. Results for the Breast Cancer Wisconsin dataset. Left: Posterior samples of the standard parameterization of \mathbf{W} (in the background) and posterior samples of \mathbf{W} parameterized by \mathbf{U} and Σ (foreground). Right: Posterior samples of our parameterization and white dots are the classical PCA solution. 4000 samples for each row of \mathbf{W} are shown in different colors. For clarity, we only show 10 randomly chosen rows instead of all 30.

direct parameterization of \mathbf{W} and the parameterization of \mathbf{W} by \mathbf{U} and Σ and again, our model is able to uniquely identify \mathbf{W} . The right plot contains the uniquely identified \mathbf{W} and the classical PCA solution (white dots). Again, our model enriches the classical PCA solution with uncertainty estimates. Samples were drawn from 4 independent chains which all converged to the same solution.

3.2. Computational Complexity

Compared to probabilistic PCA we compute the $D \times Q$ loading matrix \mathbf{W} from parameters \mathbf{v} for the principal rotation and σ^2 for the principal variances. This is achieved via a sequence of Q Householder transformations which, as detailed in Shepard et al. (2014), can be achieved in $\mathcal{O}(DQ^2)$ steps. The same complexity applies when computing gradients with respect to the model parameters. Overall, this increases the total computation time of probabilistic PCA by a factor of Q , i.e. the number of latent dimensions, compared to the standard parameterization. The scaling with respect to the number of data points is unchanged.

When sampling from the model, the adaptive NUTS sampler,

which is implemented in Stan, needs fewer leap-frog steps per sample because the rotational symmetry is removed from the posterior and does not need to be explored. Numerically, this leads to slightly lower wall clock times overall. From this perspective, our implementation is as efficient as other sampling methods proposed for probabilistic PCA or Bayesian factor analysis. If scalability becomes an issue, our model can easily be used together with more scalable approximation methods such as variational Bayes.

Pourzanjani et al. (2017) proposed an alternative parameterization which removes the rotational symmetry by means of Givens rotations. This approach requires a Jacobian adjustment to account for the change of measure in order to ensure the uniform Haar measure. This requires the computation of the absolute determinant of an $D \cdot Q$ dimensional matrix. Thus, compared to their approach we achieve a substantial speedup as, according to theorem 2, no Jacobian adjustment is required in our model.

3.3. Other than Gaussian Priors

The prior in the paper has been chosen to aid comparison with standard (classical) PCA. Using the SVD, we decompose the prior into a rotation, i.e. the principal axis, and a diagonal matrix, containing the principal components/variances. Thus, the parameters of our model are easily and well interpretable. For the Gaussian prior, which is often chosen in Bayesian factor analysis to support Gibbs sampling, we show that the rotation and variances are independent. Furthermore, the rotation is uniformly distributed according to the Haar measure.

Other priors of this structure can be constructed without issues. In particular, the interpretation of the variance parameters σ^2 as the principal components helps to this end. For instance, automatic relevance determination is readily accomplished by putting a shrinkage/sparsity prior on the principle components. Interestingly, such a prior would be very different from priors arising from imposing sparsity on the factor loadings as suggested in sparse Bayesian factor analysis (Bhattacharya & Dunson, 2011). In this case, the induced distribution on the principal components is still sparse - at least numerically, as we are not aware of any theoretical results. More importantly though, the induced rotation is not distributed uniformly anymore. Thus, such a prior gives rise to an implicit a-priori preference for certain principal axes. Note that this does not break the rotational symmetry of the model which arises from the rotation of the latent space. Instead it imposes a preference for certain directions in data space which, as we believe, is actually undesirable in many cases. We view it as a strength of our model that interpretable priors can be chosen independently for the rotation and variances.

4. Extension to Non-linear Models

4.1. Gaussian Process Latent Variable Model

Since our method works well for the linear PPCA model, the question arises, whether we could as well improve the posterior sampling of non-linear dimensionality reduction techniques. Lawrence (2005) introduced an extension of PPCA called the Gaussian Process Latent Variable Model (GP-LVM). The model starts with the dual formulation of PPCA. In this case, instead of marginalizing over \mathbf{X} in equation (2.3), we marginalize over $\mathbf{W} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The resulting likelihood then takes the form

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{d=1}^D \mathcal{N}(\mathbf{Y}_{:,d}|\boldsymbol{\mu}, \mathbf{K} + \sigma^2 \mathbf{I}), \quad (4.1)$$

where $\mathbf{K} = \mathbf{X}\mathbf{X}^T$ and $\mathbf{K}_{ij} = \mathbf{X}_{i,:}^T \mathbf{X}_{j,:}$. Identifying this with the covariance kernel of Gaussian process (Rasmussen, 2006), the model is generalized by replacing the linear kernel \mathbf{K}_{ij} with a non-linear one. In this way, a non-linear dimensionality reduction method is obtained, the GP-LVM. One of the most used covariance kernels is the exponentiated quadratic kernel (Rasmussen, 2006)

$$k_{SE}(\mathbf{x}, \mathbf{x}') = \sigma_{SE}^2 \exp(-0.5\|\mathbf{x} - \mathbf{x}'\|_2^2/l^2), \quad (4.2)$$

where σ_{SE}^2 is the kernel variance and l the lengthscale. We will use this covariance function in our analysis as well.

As suggested by Lawrence (2005), we can optimize equation (4.1) $\log p(\mathbf{Y}|\mathbf{X})$ with respect to the latent positions and the hyperparameters. Optimization can easily lead to overfitting and a fully Bayesian treatment of the model would be preferable. Therefore, we need a prior distribution $p(\mathbf{X})$ on \mathbf{X} and using Bayes rule, the posterior is given as $p(\mathbf{X}|\mathbf{Y}) = p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})/p(\mathbf{Y})$. As in the case of PPCA the posterior is not analytically tractable and, as before, we resort to sampling in order to approximate it. Note that the kernel k_{SE} only depends on the distance between points \mathbf{x} and \mathbf{x}' and thus the model likelihood is again invariant under rotations of the latent space \mathbf{X} .

4.2. Model Comparison

To test our model in the non-linear case, we again take the Breast Cancer Wisconsin dataset. For the GP-LVM the input is the transposed matrix $\mathbf{Y} \in \mathbb{R}^{N \times D}$, where $N = 30$ and $D = 569$. We standardized the data and set σ_{SE} and l to one and only sample the latent space. First, we fit the data with the standard parameterization, where we directly have the \mathbf{X} as parameters and then we reparameterize \mathbf{X} by \mathbf{U} and $\boldsymbol{\Sigma}$ and sample \mathbf{v} and $\boldsymbol{\sigma}$ as described in sections 2.4 and 2.5.

Figure 4.1 shows the results for the exponentiated quadratic kernel for $Q = 2$. As before, we used Stan to sample from

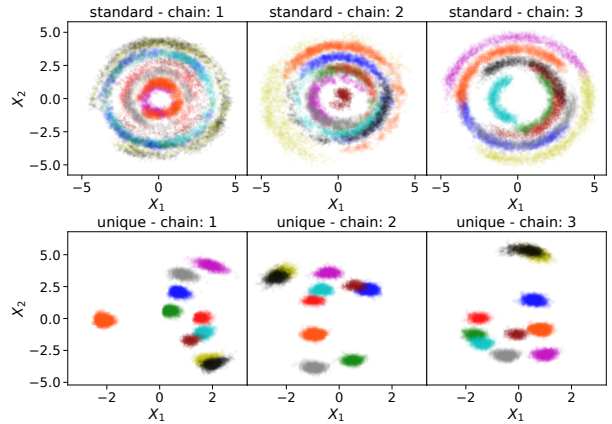


Figure 4.1. Posterior samples from for the GP-LVM. The top row shows the results for the standard parameterization for \mathbf{X} and the bottom row shows the results of our suggested parameterization. 2000 samples per chain for each row of \mathbf{X} are shown in different colors. Again, for clarity, we only show 10 randomly chosen rows instead of all 30.

the posterior on three independent chains. The top row shows the samples from the different chains for the standard parameterization. One can clearly see the rotational symmetry. The bottom row shows the samples with our suggested parameterization. As expected, there is no rotational symmetry anymore. However, due to the models complexity/flexibility many local minima arise and the chains with both (the standard and our suggested) parameterizations do not converge to the same posterior.

5. Conclusion

We suggested a new parameterization for the mapping \mathbf{W} in PPCA taking latent to observed points, which uniquely identifies the principal components even though the likelihood and the posterior (for the Bayesian case) are rotationally symmetric. We have shown how to parameterize the model via the singular vectors and values of \mathbf{W} and how to set the prior on the new parameters such that the model is not changed compared to a standard Gaussian prior on \mathbf{W} directly. Furthermore, we provided an efficient implementation via Householder transformations.

Thereafter, we tested the model on a synthetic and the Breast Cancer Wisconsin dataset. Our model was able to uniquely reconstruct the true parameters that created the synthetic data. On both datasets our model provided similar results as classical PCA, additionally containing the uncertainty of the estimates as well. In the end, we showed that our method can also be used to remove the symmetry in the latent space of a non-linear GP-LVM model. Overall, we believe that our approach, thanks to its known prior distribution and computational efficiency, can be quite useful for latent space models with rotation symmetric likelihoods.

Acknowledgements

The authors thank Dr. h.c. Maucher for funding their positions.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- Bai, Z. D., Miao, B. Q., and Pan, G. M. On asymptotics of eigenvectors of large sample covariance matrix. *Ann. Probab.*, 35(4):1532–1572, 07 2007.
- Betancourt, M. A Conceptual Introduction to Hamiltonian Monte Carlo. *ArXiv e-prints*, January 2017.
- Bhattacharya, A. and Dunson, D. B. Sparse bayesian infinite factor models. *Biometrika*, 98 2:291–306, 2011.
- Bishop, C. M. *Pattern Recognition and Machine Learning*. Springer, New York, 1st edition, 2006.
- Byrne, S. and Girolami, M. Geodesic monte carlo on embedded manifolds. *Scandinavian Journal of Statistics*, 40 (4):825–845, 2013.
- Carpenter, B., Gelman, A., Hoffman, M., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J., Li, P., and Riddell, A. Stan: A probabilistic programming language. *Journal of Statistical Software, Articles*, 76(1):1–32, 2017. ISSN 1548-7660.
- Hoff, P. D. Simulation of the matrix Bingham–von Mises–Fisher distribution, with applications to multivariate and relational data. *Journal of Computational and Graphical Statistics*, 18(2):438–456, 2009.
- James, O. and Lee, H.-N. Concise Probability Distributions of Eigenvalues of Real-Valued Wishart Matrices. *ArXiv e-prints*, February 2014.
- Jöreskog, K. G. A general approach to confirmatory maximum likelihood factor analysis. *Psychometrika*, 34(2): 183–202, Jun 1969. ISSN 1860-0980.
- Lawrence, N. Probabilistic non-linear principal component analysis with gaussian process latent variable models. *J. Mach. Learn. Res.*, 6:1783–1816, December 2005. ISSN 1532-4435.
- Mezzadri, F. How to generate random matrices from the classical compact groups. *Notices of the American Mathematical Society*, 54(5):592 – 604, 5 2007. ISSN 0002-9920.
- Millsap, R. E. When trivial constraints are not trivial: The choice of uniqueness constraints in confirmatory factor analysis. *Structural Equation Modeling: A Multidisciplinary Journal*, 8(1):1–17, 2001.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Peeters, C. F. W. Rotational uniqueness conditions under oblique factor correlation metric. *Psychometrika*, 77(2): 288–292, Apr 2012. ISSN 1860-0980.
- Pourzanjani, A. A., Jiang, R. M., Mitchell, B., Atzberger, P. J., and Petzold, L. R. General Bayesian Inference over the Stiefel Manifold via the Givens Transform. *ArXiv e-prints*, October 2017.
- Rasmussen, C. E. *Gaussian processes for machine learning*. MIT Press, 2006.
- Shepard, R., Gidofalvi, G., and Brozell, S. R. The multifacet graphically contracted function method. II. a general procedure for the parameterization of orthogonal matrices and its application to arc factors. *The Journal of Chemical Physics*, 141(6):064106, 2014.
- Stan Development Team. Stan modeling language users guide and reference manual, version 2.18.0, 2018.
- Tipping, M. E. and Bishop, C. M. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, Series B*, 61:611–622, 1999.
- Uhlig, H. On singular wishart and singular multivariate beta distributions. *Ann. Statist.*, 22(1):395–405, 03 1994.
- Šmídl, V. and Quinn, A. On bayesian principal component analysis. *Computational Statistics & Data Analysis*, 51 (9):4101–4123, 2007.