

Rotation Invariant Householder Parameterization for Bayesian PCA

Rajbir-Singh Nirwan, Nils Bertschinger

June 11, 2019



Outline

- Probabilistic PCA (PPCA)
- Non-identifiability issue of PPCA
- Conceptual solution to the problem
- Implementation
- Results

Probabilistic PCA

- Classical PCA

Formulated as a projection from data space \mathbf{Y} to a lower dimensional latent space \mathbf{X}

$$\mathbf{Y} \in \mathbb{R}^{N \times D} \rightarrow \mathbf{X} \in \mathbb{R}^{N \times Q}$$

Latent space: maximizes variance of projected data, minimizes MSE

Probabilistic PCA

- Classical PCA

Formulated as a projection from data space \mathbf{Y} to a lower dimensional latent space \mathbf{X}

$$\mathbf{Y} \in \mathbb{R}^{N \times D} \quad \rightarrow \quad \mathbf{X} \in \mathbb{R}^{N \times Q}$$

Latent space: maximizes variance of projected data, minimizes MSE

- Probabilistic PCA (PPCA)

Viewed as a generative model, that maps the latent space \mathbf{X} to the data space \mathbf{Y}

$$\mathbf{X} \in \mathbb{R}^{N \times Q} \quad \rightarrow \quad \mathbf{Y} \in \mathbb{R}^{N \times D}$$

$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{Y} | \mathbf{W}) = \prod_{n=1}^N \mathcal{N}(\mathbf{Y}_{n,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$$

$$\mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T \quad \forall \mathbf{R}\mathbf{R}^T = \mathbf{I}$$

Probabilistic PCA

- Classical PCA

Formulated as a projection from data space \mathbf{Y} to a lower dimensional latent space \mathbf{X}

$$\mathbf{Y} \in \mathbb{R}^{N \times D} \rightarrow \mathbf{X} \in \mathbb{R}^{N \times Q}$$

Latent space: maximizes variance of projected data, minimizes MSE

- Probabilistic PCA (PPCA)

Viewed as a generative model, that maps the latent space \mathbf{X} to the data space \mathbf{Y}

$$\mathbf{X} \in \mathbb{R}^{N \times Q} \rightarrow \mathbf{Y} \in \mathbb{R}^{N \times D}$$

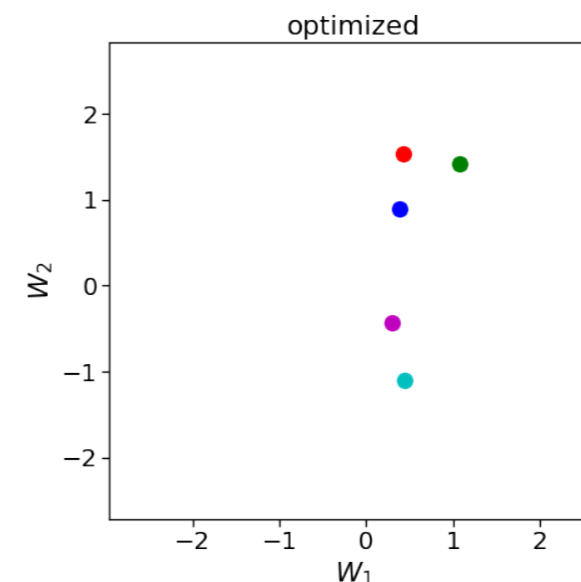
$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{Y} | \mathbf{W}) = \prod_{n=1}^N \mathcal{N}(\mathbf{Y}_{n,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$$

$$\mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T \quad \forall \mathbf{R}\mathbf{R}^T = \mathbf{I}$$

- Optimization for $D=5, Q=2$



Probabilistic PCA

- Classical PCA

Formulated as a projection from data space \mathbf{Y} to a lower dimensional latent space \mathbf{X}

$$\mathbf{Y} \in \mathbb{R}^{N \times D} \rightarrow \mathbf{X} \in \mathbb{R}^{N \times Q}$$

Latent space: maximizes variance of projected data, minimizes MSE

- Probabilistic PCA (PPCA)

Viewed as a generative model, that maps the latent space \mathbf{X} to the data space \mathbf{Y}

$$\mathbf{X} \in \mathbb{R}^{N \times Q} \rightarrow \mathbf{Y} \in \mathbb{R}^{N \times D}$$

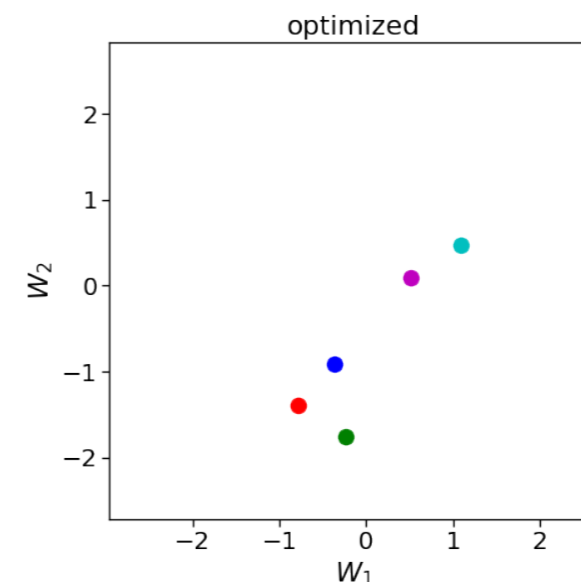
$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{Y} | \mathbf{W}) = \prod_{n=1}^N \mathcal{N}(\mathbf{Y}_{n,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$$

$$\mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T \quad \forall \mathbf{R}\mathbf{R}^T = \mathbf{I}$$

- Optimization for $D=5, Q=2$



Probabilistic PCA

- Classical PCA

Formulated as a projection from data space \mathbf{Y} to a lower dimensional latent space \mathbf{X}

$$\mathbf{Y} \in \mathbb{R}^{N \times D} \rightarrow \mathbf{X} \in \mathbb{R}^{N \times Q}$$

Latent space: maximizes variance of projected data, minimizes MSE

- Probabilistic PCA (PPCA)

Viewed as a generative model, that maps the latent space \mathbf{X} to the data space \mathbf{Y}

$$\mathbf{X} \in \mathbb{R}^{N \times Q} \rightarrow \mathbf{Y} \in \mathbb{R}^{N \times D}$$

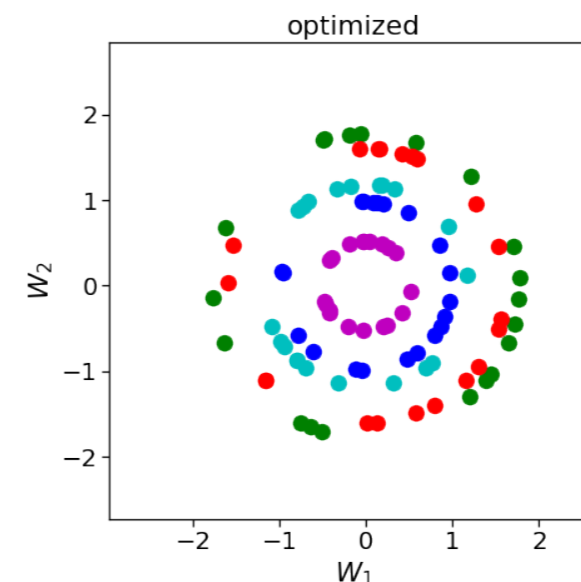
$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{Y} | \mathbf{W}) = \prod_{n=1}^N \mathcal{N}(\mathbf{Y}_{n,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$$

$$\mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T \quad \forall \mathbf{R}\mathbf{R}^T = \mathbf{I}$$

- Optimization for $D=5, Q=2$



Probabilistic PCA

- Classical PCA

Formulated as a projection from data space \mathbf{Y} to a lower dimensional latent space \mathbf{X}

$$\mathbf{Y} \in \mathbb{R}^{N \times D} \rightarrow \mathbf{X} \in \mathbb{R}^{N \times Q}$$

Latent space: maximizes variance of projected data, minimizes MSE

- Probabilistic PCA (PPCA)

Viewed as a generative model, that maps the latent space \mathbf{X} to the data space \mathbf{Y}

$$\mathbf{X} \in \mathbb{R}^{N \times Q} \rightarrow \mathbf{Y} \in \mathbb{R}^{N \times D}$$

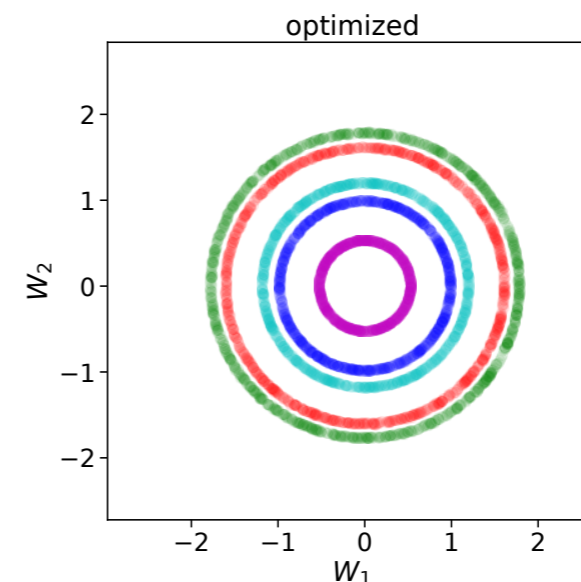
$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{Y} | \mathbf{W}) = \prod_{n=1}^N \mathcal{N}(Y_{n,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$$

$$\mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T \quad \forall \mathbf{R}\mathbf{R}^T = \mathbf{I}$$

- Rotation invariant likelihood



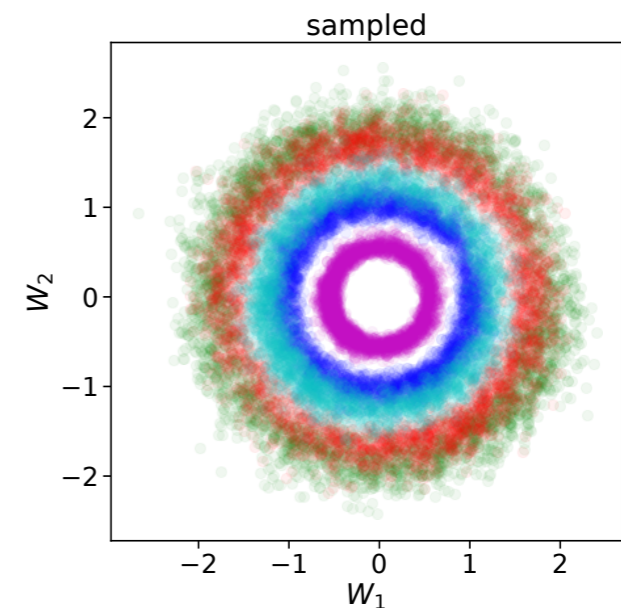
Bayesian approach to PPCA

$$p(W|Y) = \frac{p(Y|W)p(W)}{p(Y)}$$

- If prior does not break the symmetry, posterior will be rotation invariant as well
- Sampling will be challenging, posterior averages are meaningless and the interpretation of the latent space is almost impossible

Bayesian approach to PPCA

$$p(W|Y) = \frac{p(Y|W)p(W)}{p(Y)}$$



- If prior does not break the symmetry, posterior will be rotation invariant as well
- Sampling will be challenging, posterior averages are meaningless and the interpretation of the latent space is almost impossible

Solution

- Find different parameterization of the model, such that the probabilistic model is not changed

Outline of procedure

- SVD of W $WW^T = U\Sigma V^T (U\Sigma V^T)^T = U\Sigma^2 U^T$
- Fix coordinate system $V = I$
- Specify correct prior $p(U, \Sigma)$
- Sample from $p(U, \Sigma | Y)$

Solution

- Find different parameterization of the model, such that the probabilistic model is not changed

Outline of procedure

- SVD of W $WW^T = U\Sigma V^T (U\Sigma V^T)^T = U\Sigma^2 U^T$
- Fix coordinate system $V = I$
- Specify correct prior $p(U, \Sigma)$
- Sample from $p(U, \Sigma | Y)$

$W \sim \mathcal{N}(\mathbf{0}, I) \rightarrow WW^T$ is Wishart distributed

$U \sim ?$
 $\Sigma \sim ? \rightarrow U\Sigma\Sigma^T U^T$ is Wishart distributed

$$\begin{matrix} U \sim ? \\ \Sigma \sim ? \end{matrix} \rightarrow U \Sigma \Sigma^T U^T \text{ Wishart}$$

Theory

- Since $\mathbf{U}, \mathbf{\Sigma}$ is SVD of \mathbf{W} and $\mathbf{U}, \mathbf{\Sigma}^2$ is eigenvalue decomposition of $\mathbf{W}\mathbf{W}^T \rightarrow \mathbf{U}$ is eigenvector matrix

$$U \in \mathcal{V}_{Q,D} \quad \text{Stiefel manifold} \quad \mathcal{V}_{Q,D} = \{U \in \mathbb{R}^{D \times Q} \mid U^T U = I\}$$

Eigenvectors of Wishart matrix are distributed uniformly in space of orthogonal matrices (Blai (2007), Uhlig (1994))

$\rightarrow \mathbf{U}$ is uniformly distributed on the Stiefel manifold

$$\begin{matrix} U \sim ? \\ \Sigma \sim ? \end{matrix} \rightarrow U \Sigma \Sigma^T U^T \text{ wishart}$$

Theory

- Since \mathbf{U}, Σ is SVD of \mathbf{W} and \mathbf{U}, Σ^2 is eigenvalue decomposition of $\mathbf{W}\mathbf{W}^T \rightarrow \mathbf{U}$ is eigenvector matrix

$$U \in \mathcal{V}_{Q,D} \quad \text{Stiefel manifold} \quad \mathcal{V}_{Q,D} = \{U \in \mathbb{R}^{D \times Q} \mid U^T U = I\}$$

Eigenvectors of Wishart matrix are distributed uniformly in space of orthogonal matrices (Blai (2007), Uhlig (1994))

$\rightarrow \mathbf{U}$ is uniformly distributed on the Stiefel manifold

- Square of ordered eigenvalue matrix Σ is distributed as (James & Lee (2014))

$$p(\lambda) = c e^{-\frac{1}{2} \sum_{q=1}^Q \lambda_q} \prod_{q=1}^Q \left(\lambda_q^{\frac{D-Q-1}{2}} \prod_{q'=q+1}^Q |\lambda_q - \lambda_{q'}| \right)$$

$$p(\sigma_1, \dots, \sigma_Q) = c e^{-\frac{1}{2} \sum_{q=1}^Q \sigma_q^2} \prod_{q=1}^Q \left(\sigma_q^{D-Q-1} \prod_{q'=q+1}^Q |\sigma_q^2 - \sigma_{q'}^2| \right) \prod_{q=1}^Q 2\sigma_q$$

Implementation

- Need: $U \sim$ uniform on Stiefel $\mathcal{V}_{Q,D}$
 $\Sigma \sim p(\Sigma) \leftarrow$ easy, since we know the analytic exp for density

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 $O(D)$.

Mezzadri (2007)

How to uniformly sample U on $\mathcal{V}_{Q,D}$

for $n = D : 1$

$\mathbf{v}_n \sim$ uniform on \mathbb{S}^{n-1}

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

$$\tilde{\mathbf{H}}_n(\mathbf{v}_n) = -\text{sgn}(\mathbf{v}_{n1}) (\mathbf{I} - 2\mathbf{u}_n\mathbf{u}_n^T)$$

$$\mathbf{H}_n = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{H}}_n \end{pmatrix}$$

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

Implementation

The full generative model for Bayesian PPCA:

$$\mathbf{v}_D, \dots, \mathbf{v}_{D-Q+1} \sim \mathcal{N}(0, \mathbf{I})$$

$$\sigma \sim p(\sigma)$$

$$\boldsymbol{\mu} \sim p(\boldsymbol{\mu})$$

$$\mathbf{U} = \prod_{q=1}^Q \mathbf{H}_{D-q+1} \left(\mathbf{v}_{D-q+1} \right)$$

$$\boldsymbol{\Sigma} = \text{diag}(\sigma)$$

$$\mathbf{W} = \mathbf{U}\boldsymbol{\Sigma}$$

$$\sigma \text{ noise} \sim p \left(\sigma \text{ noise} \right)$$

$$\mathbf{Y} \sim \prod_{n=1}^N \mathcal{N} \left(\mathbf{Y}_{n,:} \mid \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \text{ noise } \mathbf{I} \right)$$

Results

Synthetic Dataset

- Construction

$$(N, D, Q) = (150, 5, 2)$$

$$X \sim \mathcal{N}(\mathbf{0}, I) \in \mathbb{R}^{N \times Q}$$

$$U \sim \text{uniform on Stiefel } \mathcal{V}_{Q,D}$$

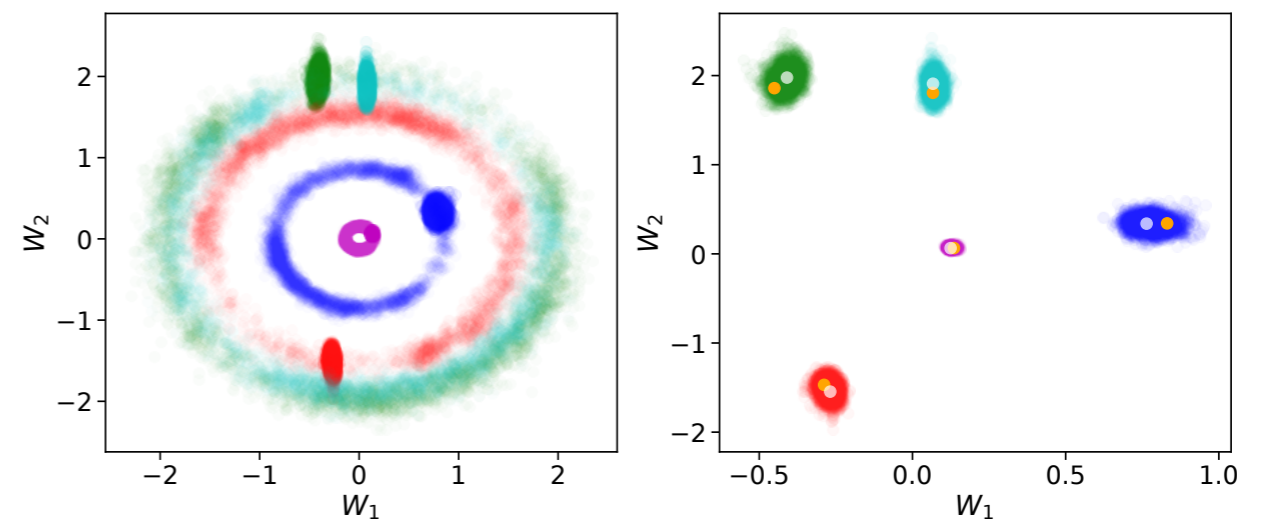
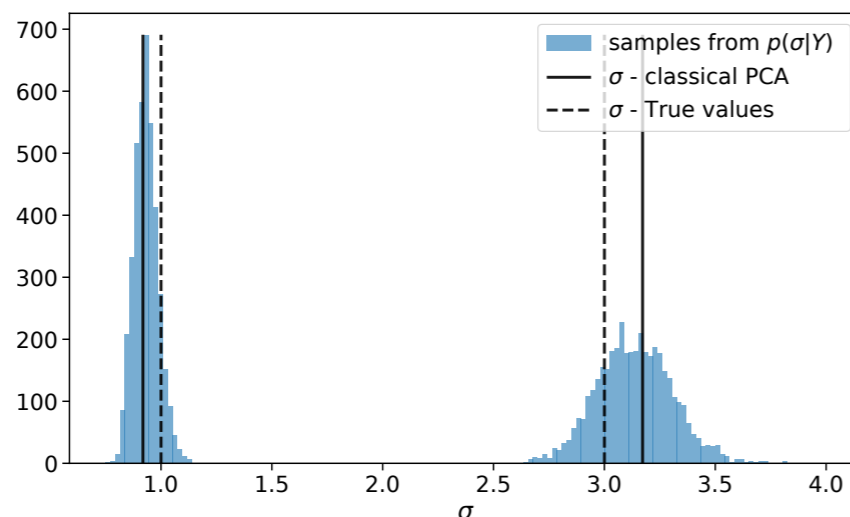
$$\epsilon \sim \mathcal{N}(0, 0.01) \in \mathbb{R}^{N \times D}$$

$$\Sigma = \text{diag}(\sigma_1, \sigma_2) = \text{diag}(3.0, 1.0)$$

$$W = U\Sigma \in \mathbb{R}^{D \times Q}$$

$$Y = XW^T + \epsilon$$

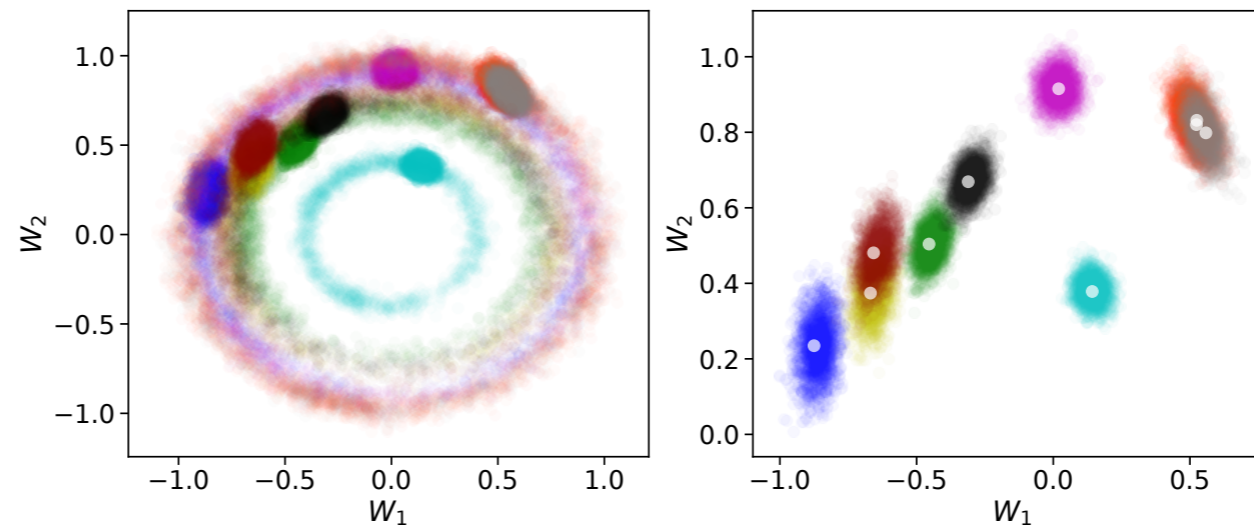
- Inference



Results

Breast Cancer Wisconsin Dataset $(N, D) = (569, 30)$

- Bayesian PCA



- Advantages

- Breaks the rotation symmetry without changing the probabilistic model
- Enrichment of the classical PCA solution with uncertainty estimates
- Decomposition of prior into rotation and principle variances
 - Allows to construct other priors without issues
 - Sparsity prior on principle variances without a-priori rotation preference
 - If desired a-priori rotation preference without affecting the variances

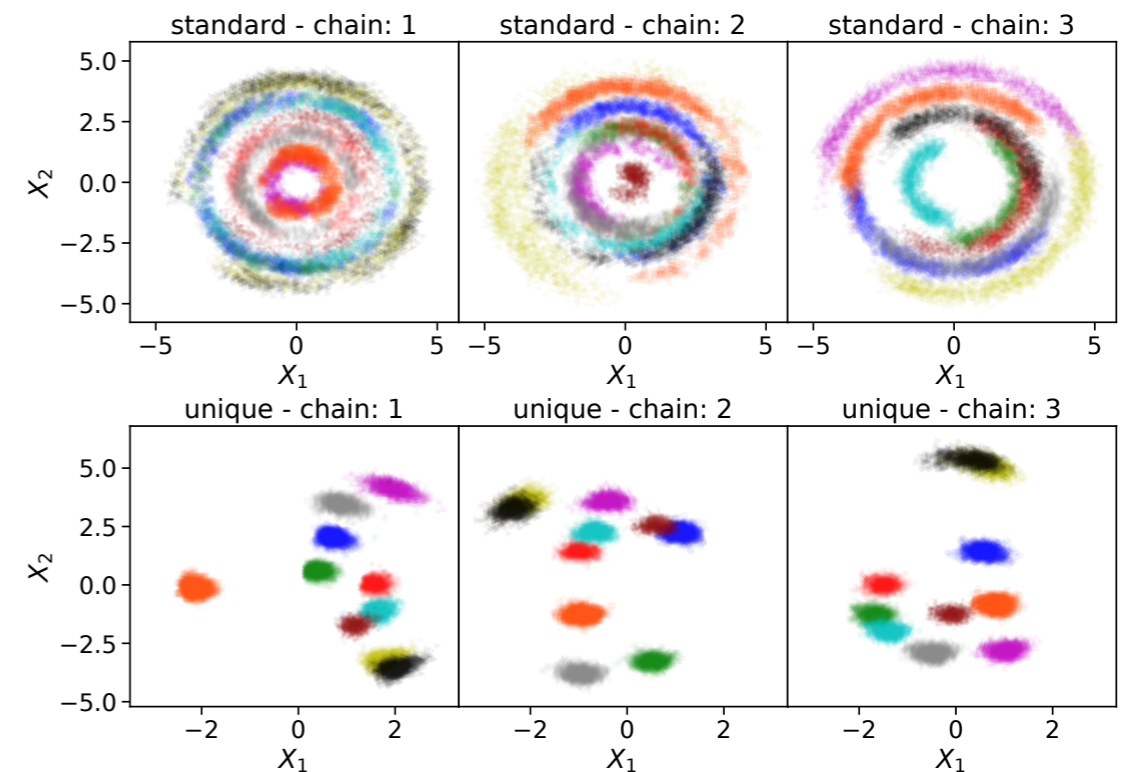
Extension to non-linear models

- GPLVM with the same rotation invariant problem

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

$$\mathbf{K} = \mathbf{X}\mathbf{X}^T, \quad \mathbf{K}_{ij} = \mathbf{X}_{i,:}^T \mathbf{X}_{j,:} = k(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})$$

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



- No rotation symmetry in the posterior for the suggested parameterization
- Different chains converge to different solutions due to increased model complexity

Conclusion

- Suggested new parameterization for \mathbf{W} in PPCA, which uniquely identifies principle components even though the likelihood and the posterior are rotationally symmetric
- Showed 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}
- Provided an efficient implementation via Householder transformations (no Jacobian correction needed)
- New parameterization allows for other interpretable priors on rotation and principle variances
- Extended to non-linear models and successfully solved the rotation problem there as well

**Thanks for your
attention!**

Supervisor: Prof. Dr. Nils Bertschinger

Funder: Dr. h. c. Helmut O. Maucher