

# Rotation Invariant Householder Parameterization for Bayesian PCA

Rajbir-Singh Nirwan  
Sept 04, 2019



# Outline

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

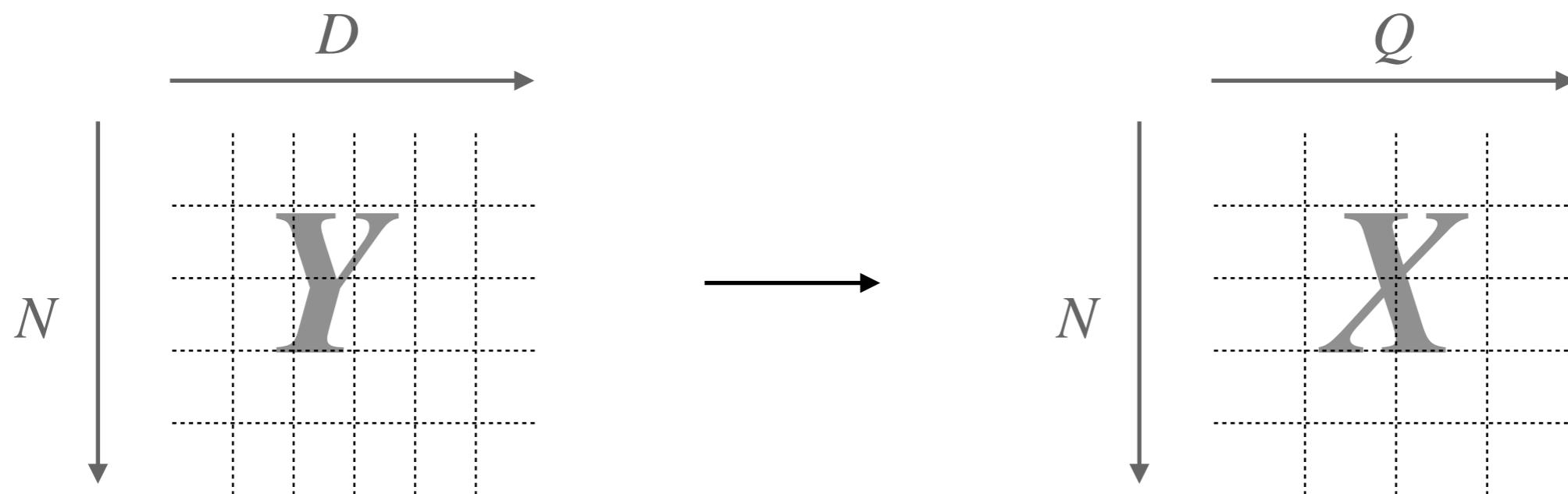
# Probabilistic PCA

- Classical PCA

Formulated as a projection from data space  $Y$  to a lower dimensional latent space  $X$

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

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  $\mathbf{X}$

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

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  $\mathbf{Y}$

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

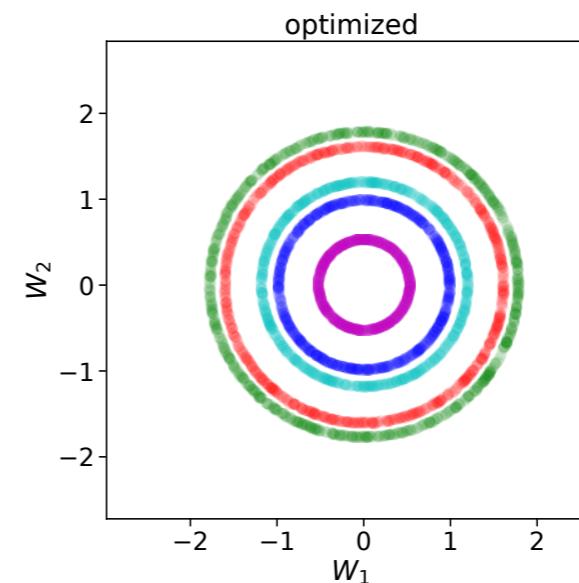
- Rotation invariant likelihood

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

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

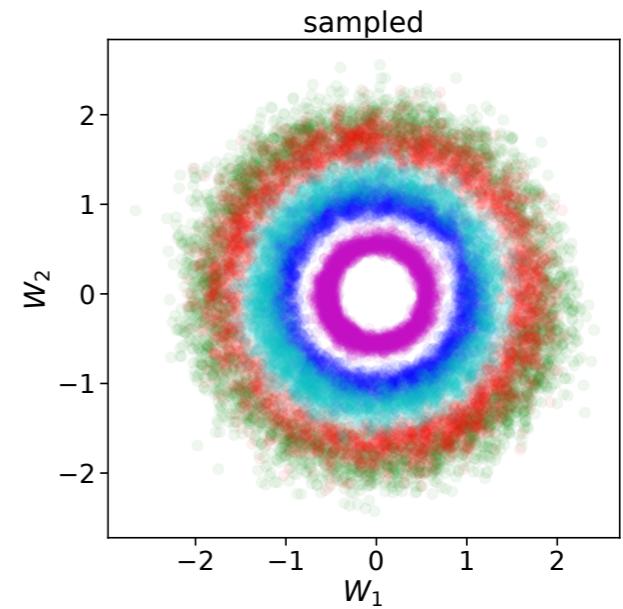
$$p(\mathbf{Y} | \mathbf{W}) = \prod_{n=1}^N \mathcal{N}(Y_{n,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$$

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



# 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  $\mathbf{W}$  
$$\mathbf{WW}^T = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T (\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T)^T = \mathbf{U}\boldsymbol{\Sigma}^2\mathbf{U}^T$$
- Fix coordinate system 
$$\mathbf{V} = \mathbf{I}$$
- Specify correct prior 
$$p(\mathbf{U}, \boldsymbol{\Sigma})$$
- Sample from 
$$p(\mathbf{U}, \boldsymbol{\Sigma} | \mathbf{Y})$$

---

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

$\mathbf{U} \sim ?$      $\rightarrow \mathbf{U}\boldsymbol{\Sigma}\boldsymbol{\Sigma}^T\mathbf{U}^T$  is **Wishart distributed**

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

# Theory

- Since  $U, \Sigma$  is SVD of  $W \rightarrow U$  is a orthogonal matrix

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

- $U, \Sigma^2$  is eigenvalue decomposition of  $WW^T$

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

**$\rightarrow U$  is uniformly distributed on the Stiefel manifold**

- Distribution of the ordered eigenvalue matrix  $\Sigma^2$  of  $WW^T$  is (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:  $U \sim \text{uniform on Stiefel } \mathcal{V}_{Q,D}$   
 $\Sigma \sim p(\Sigma) \leftarrow \text{easy, since we know the analytic exp for density}$

**Theorem 2** Let  $v_D, v_{D-1}, \dots, v_1$  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  $H_n(v_n)$  be the  $n$ -th Householder transformation as defined in equation (2.20). The product

$$Q = H_D(v_D)H_{D-1}(v_{D-1})\dots H_1(v_1) \quad (2.21)$$

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  $n = D : 1$

$$v_n \sim \text{uniform on } \mathbb{S}^{n-1}$$

$$u_n = \frac{v_n + \text{sgn}(v_{n1}) \| v_n \| e_1}{\| v_n + \text{sgn}(v_{n1}) \| v_n \| e_1 \|}$$

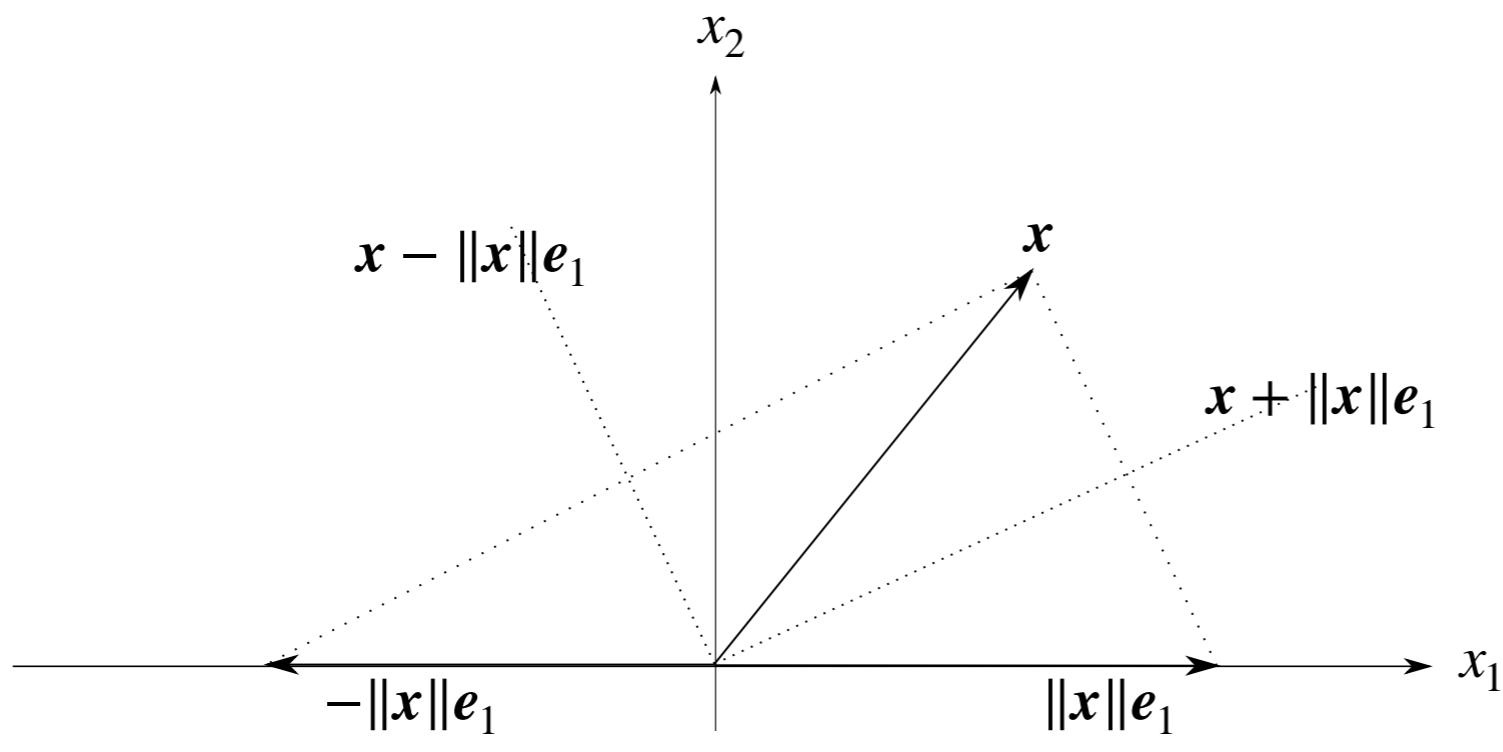
$$\tilde{H}_n(v_n) = -\text{sgn}(v_{n1})(I - 2u_n u_n^T)$$

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

$$U = H_D(v_D)H_{D-1}(v_{D-1})\dots H_1(v_1)$$

# Householder Transformations

- Used to reflect a vector in such a way that all coordinates but one disappear, e.g.: QR-decomposition



$$u = x \pm \|x\|e_1$$

$$H = \mathbf{1} - 2\hat{u}\hat{u}^T$$

$$Hx = \|x\|e_1$$

# Householder Transformations

Example for  $D = 2$

$$\boldsymbol{v}_1 = \begin{pmatrix} 0 \\ v_{11} \end{pmatrix} \quad \boldsymbol{v}_2 = \begin{pmatrix} v_{21} \\ v_{22} \end{pmatrix}$$

Construction of  $\mathbf{H}_1$

$$v_{11} \in \{-1, 1\}$$

$$\tilde{\mathbf{H}}_1 = v_{11}$$

$$\mathbf{H}_1 = \begin{pmatrix} 1 & 0 \\ 0 & v_{11} \end{pmatrix}$$

Construction of  $\mathbf{H}_2$

$$\boldsymbol{v}_2 \in \mathbb{S}^1$$

$$\boldsymbol{v}_1 = \begin{pmatrix} 0 \\ v_{11} \end{pmatrix} \quad \boldsymbol{v}_2 = \begin{pmatrix} v_{21} \\ v_{22} \end{pmatrix}$$

$$U = \mathbf{H}_2 \mathbf{H}_1 = \mathbf{H}_2 \begin{pmatrix} 1 & 0 \\ 0 & v_{11} \end{pmatrix} = \begin{pmatrix} \boldsymbol{v}_2 & \mathbf{H}_2 \boldsymbol{v}_1 \end{pmatrix}$$

for  $n = D : 1$

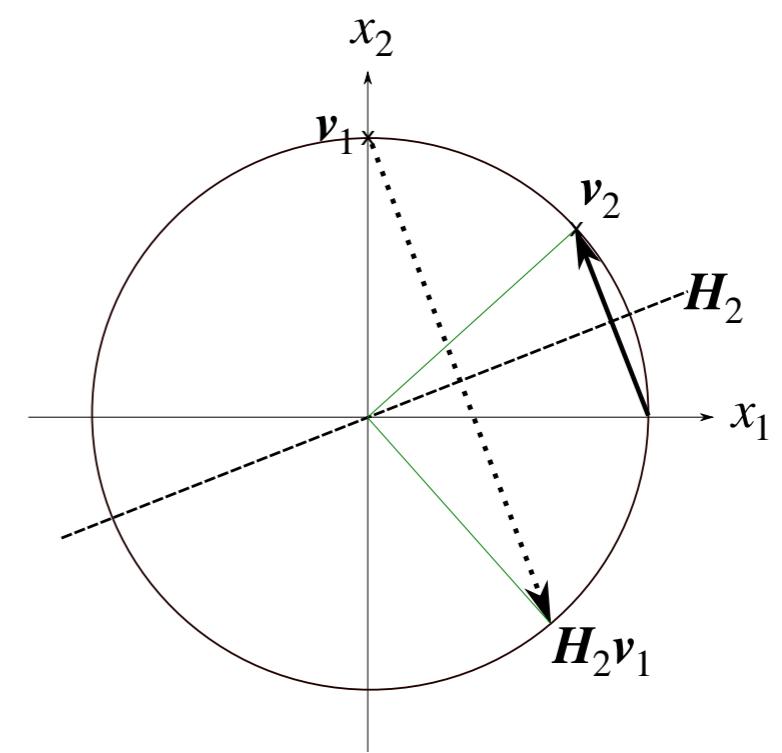
$$\boldsymbol{v}_n \sim \text{uniform on } \mathbb{S}^{n-1}$$

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

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

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

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



# Implementation

The full generative model for Bayesian PPCA:

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

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

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

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

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

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

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

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

# Results

## Synthetic Dataset

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

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

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

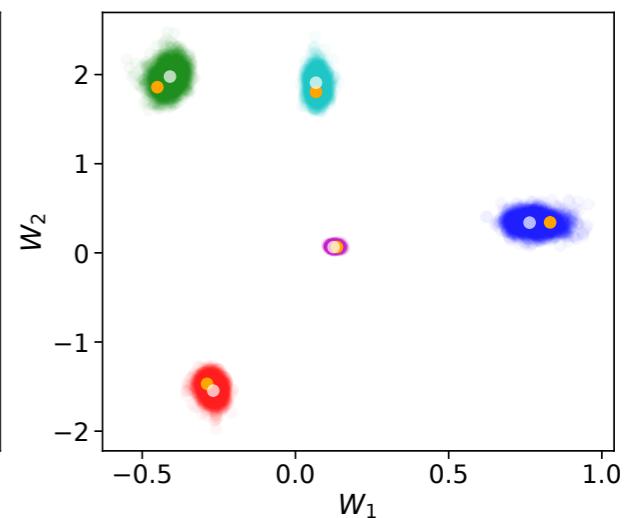
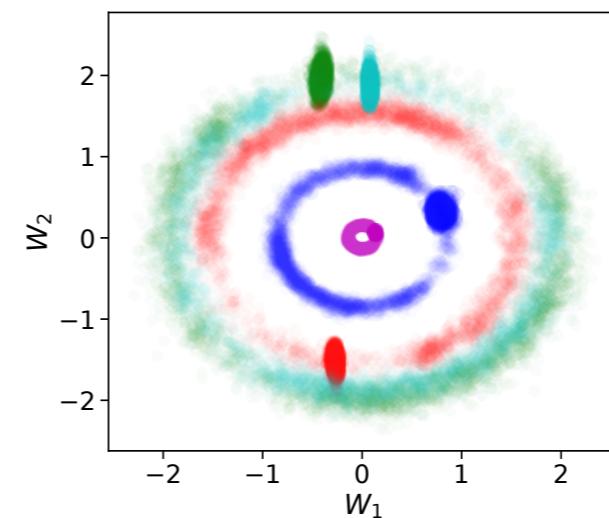
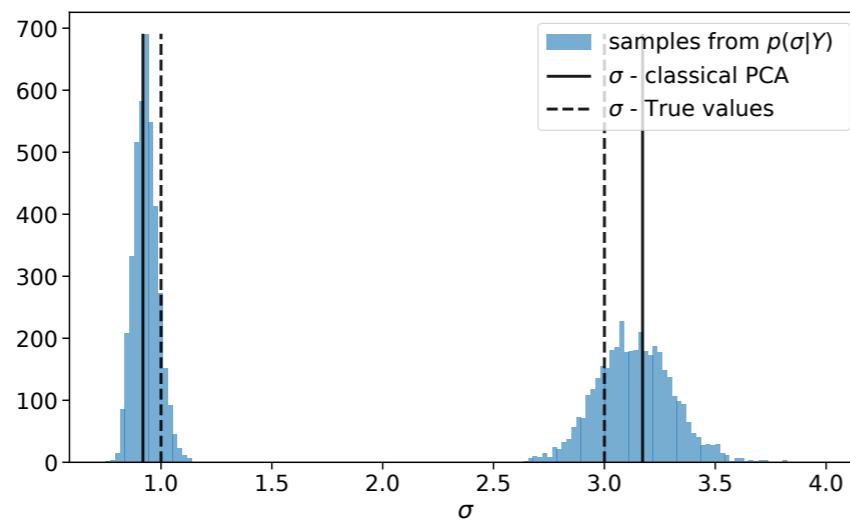
$$\epsilon \sim \mathcal{N}(0, 0.01) \in \mathbb{R}^{N \times D}$$

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

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

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

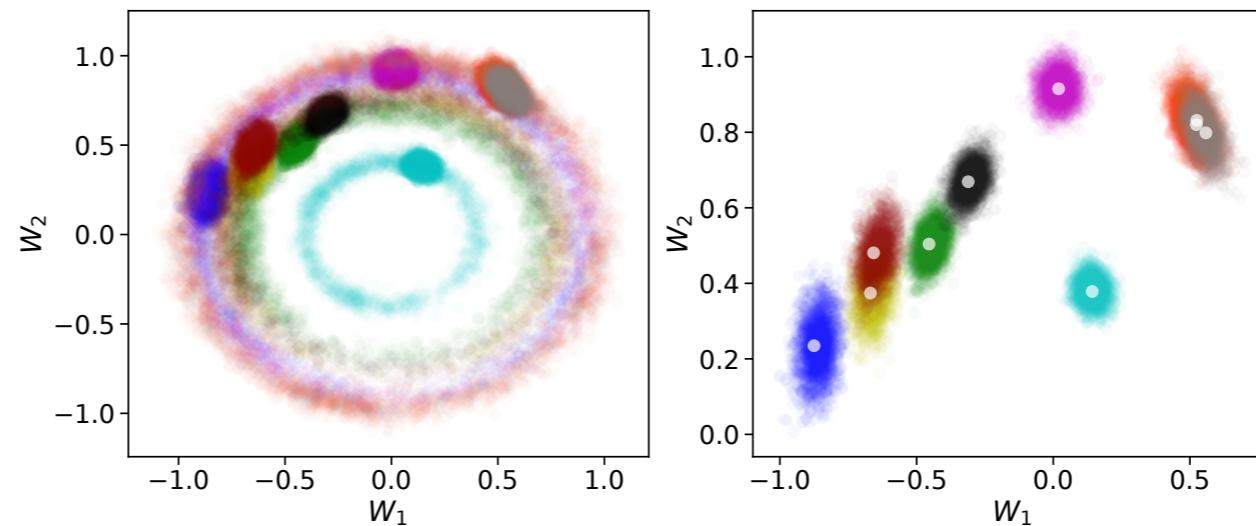
- Inference



# 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

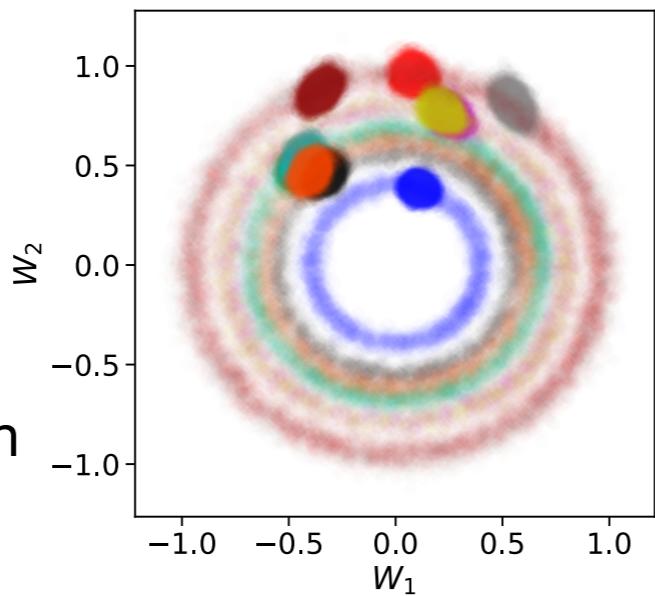
# Results

Breast Cancer

$(N, D) = (569, 30)$

Time House: 9.5 min

Time Standard: 25.6 min

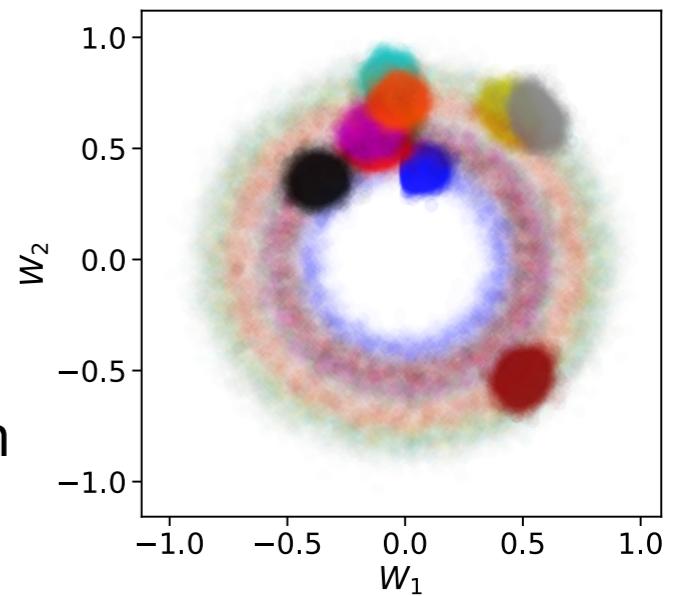


Diabetes

$(N, D) = (442, 10)$

Time House: 8.9 min

Time Standard: 2.4 min

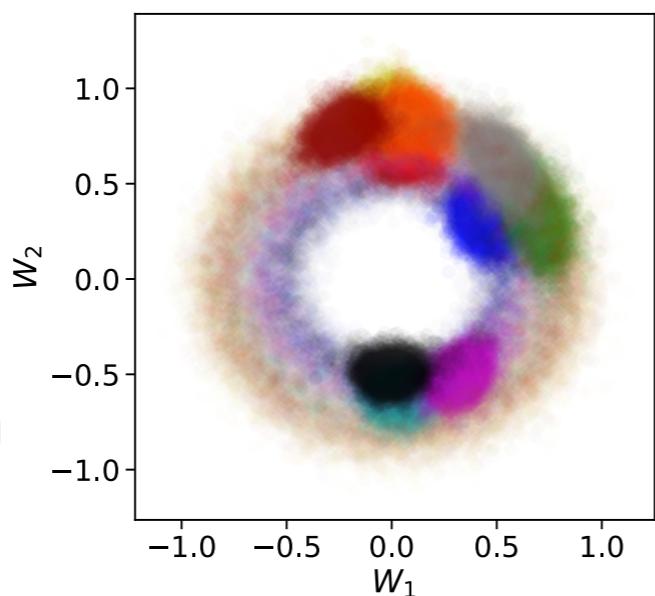


Wine

$(N, D) = (178, 13)$

Time House: 0.4 min

Time Standard: 1.2 min

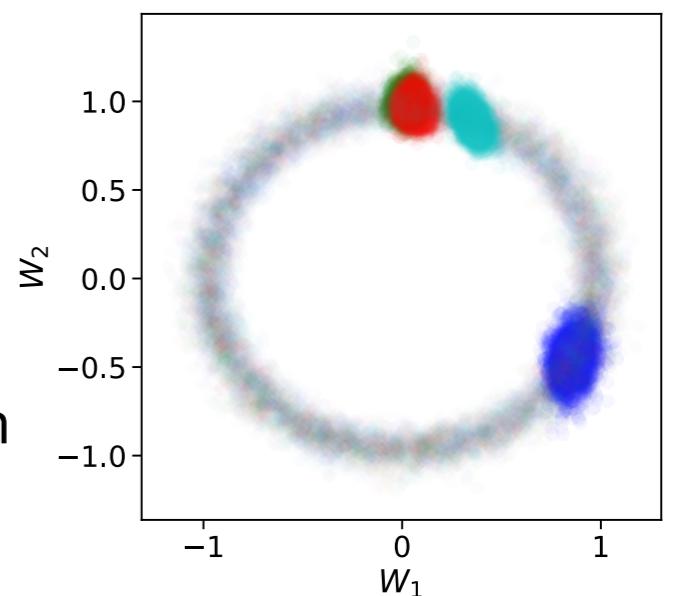


Iris

$(N, D) = (150, 4)$

Time House: 0.2 min

Time Standard: 0.8 min



# Results

**for**  $n = D : 1$

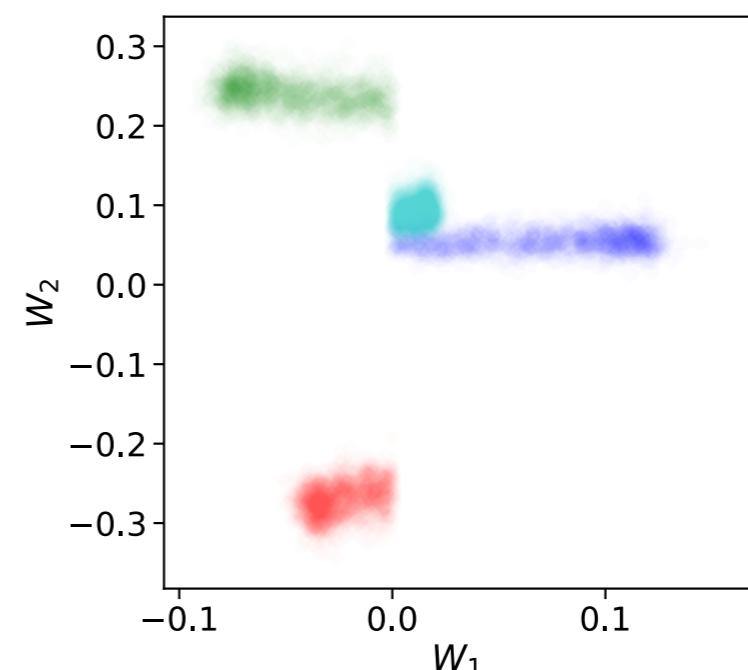
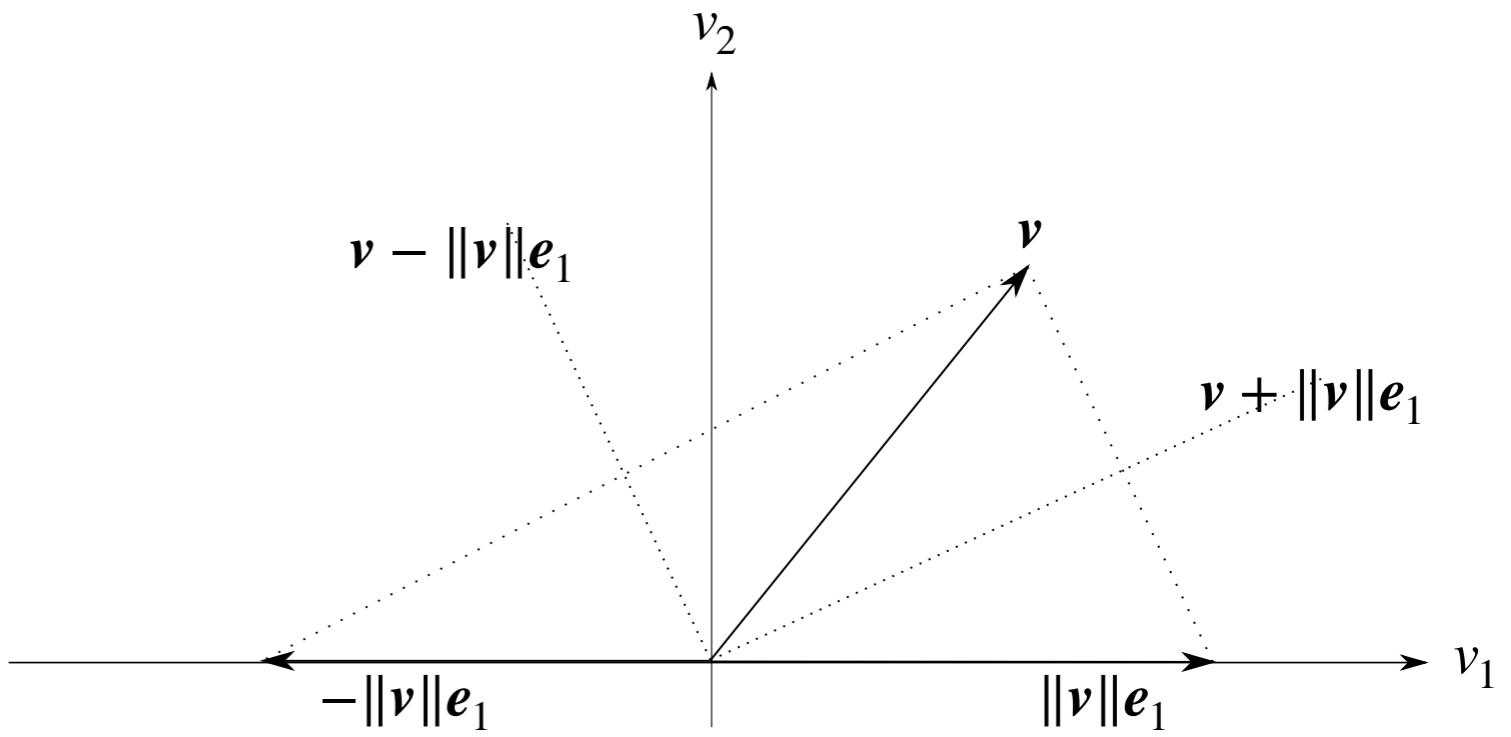
$v_n \sim \text{uniform on } \mathbb{S}^{n-1}$

$$u_n = \frac{v_n + \text{sgn}(v_{n1}) \|v_n\| e_1}{\|v_n + \text{sgn}(v_{n1}) \|v_n\| e_1\|}$$

$$\tilde{H}_n(v_n) = -\text{sgn}(v_{n1}) (\mathbf{I} - 2u_n u_n^T)$$

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

$$U = H_D(v_D) H_{D-1}(v_{D-1}) \dots H_1(v_1)$$



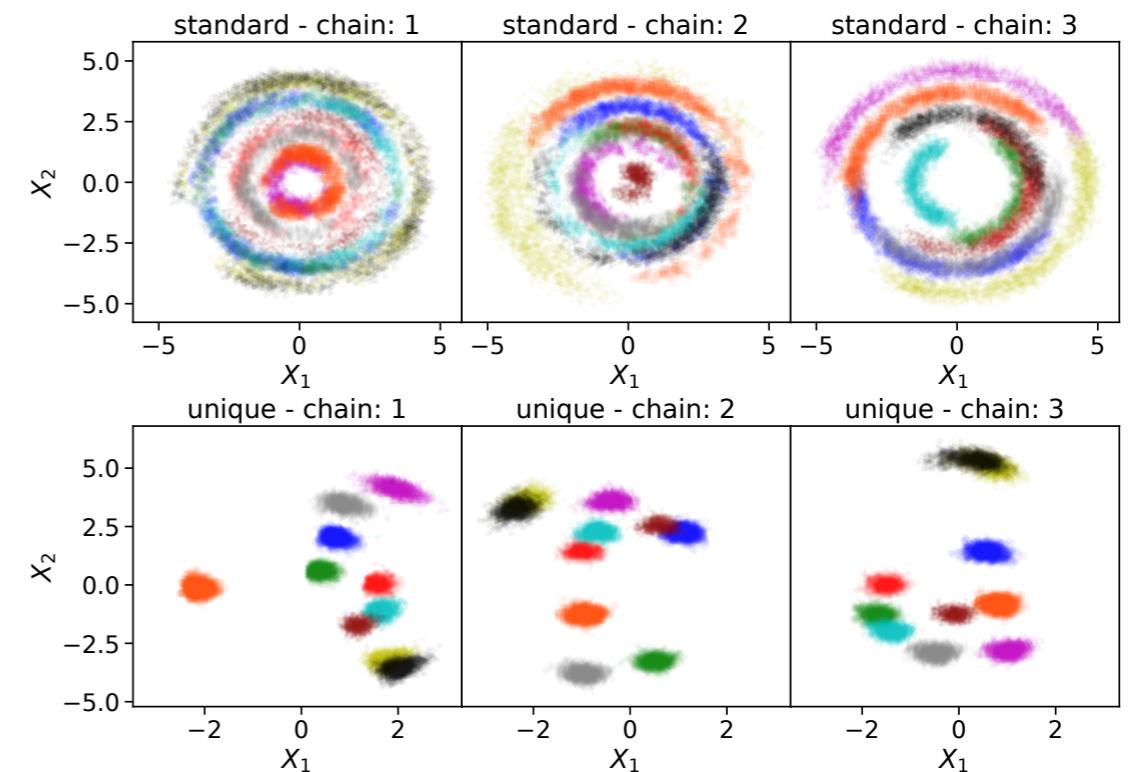
# Extension to non-linear models

- GPLVM with the same rotation invariant problem

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

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

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



- 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 is 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**