

Latent Variable Models

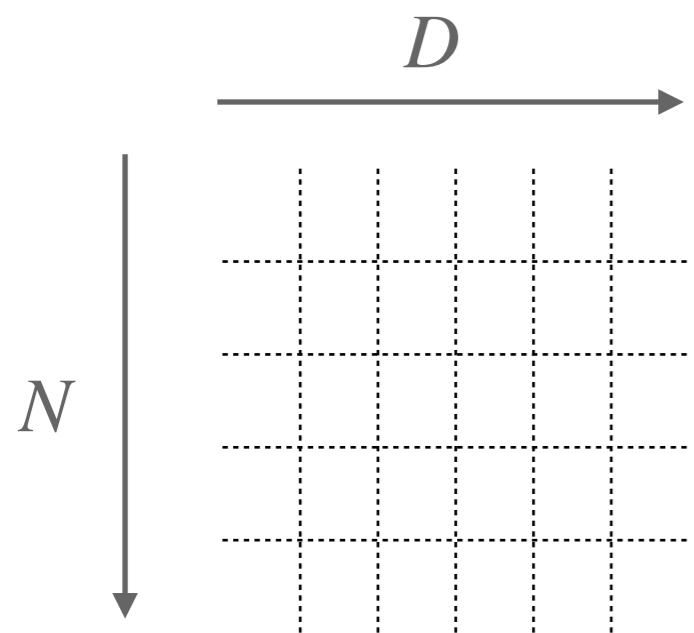
Rajbir-Singh Nirwan
Aug 27, 2019

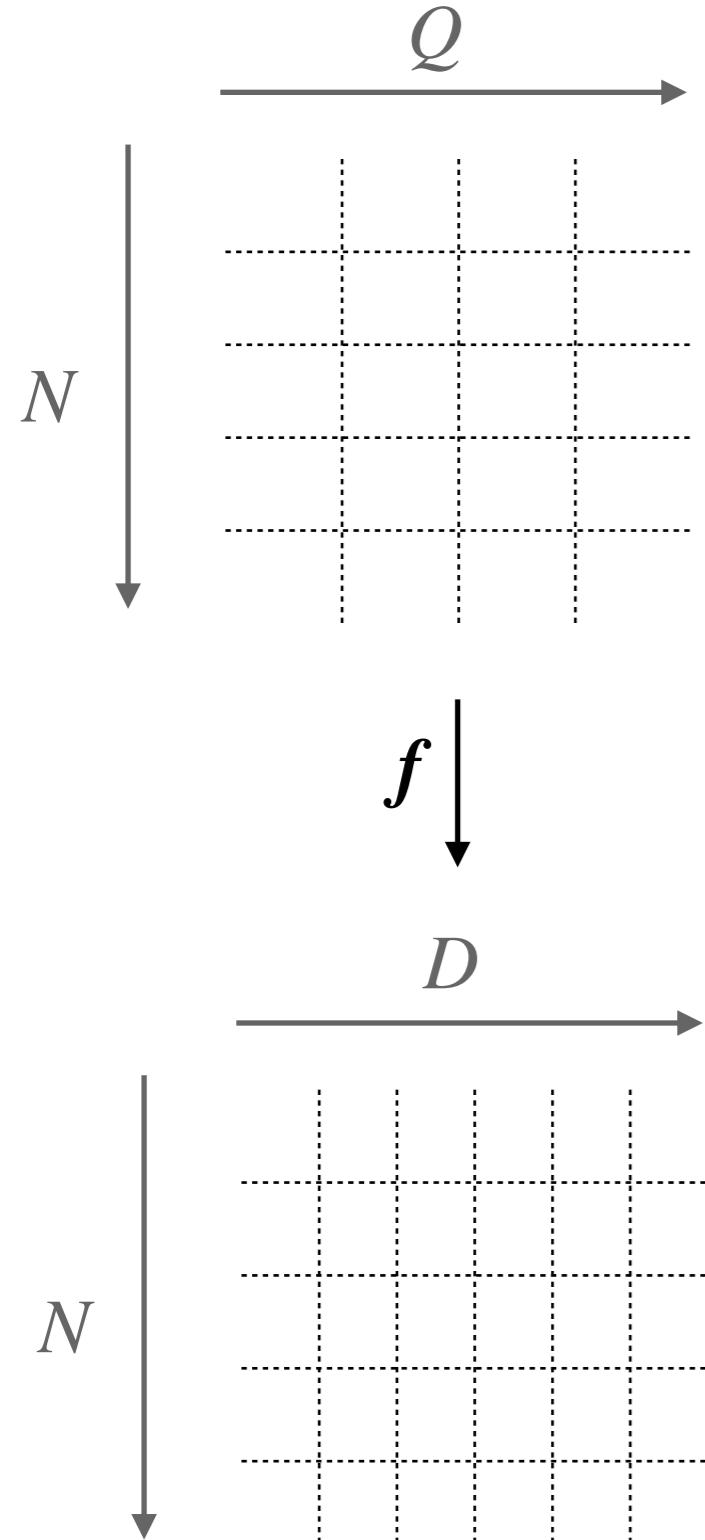
Outline

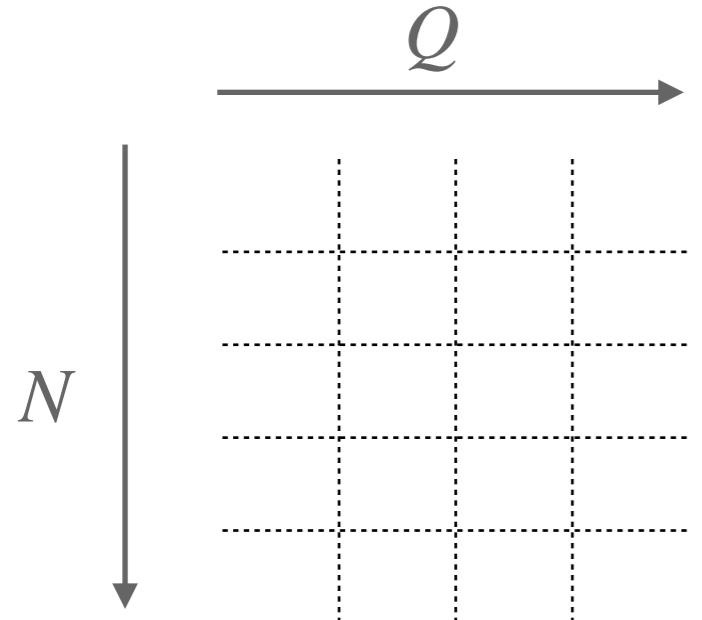
- Latent Variable Models
- Examples: PPCA, VAE
- Non-identifiability issue of PPCA
- Solution
- Summary



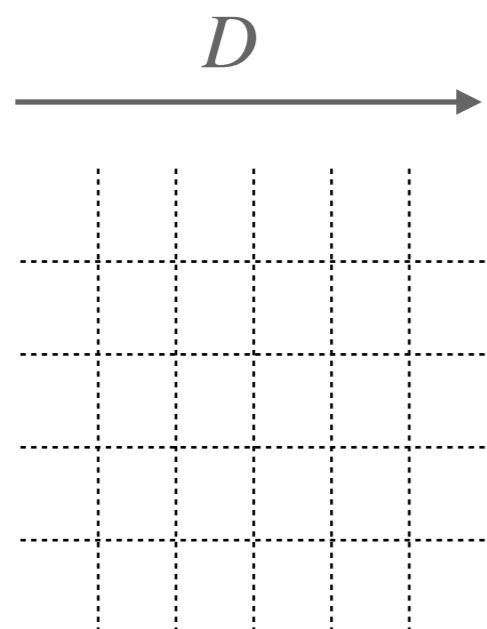
www.123rf.com/



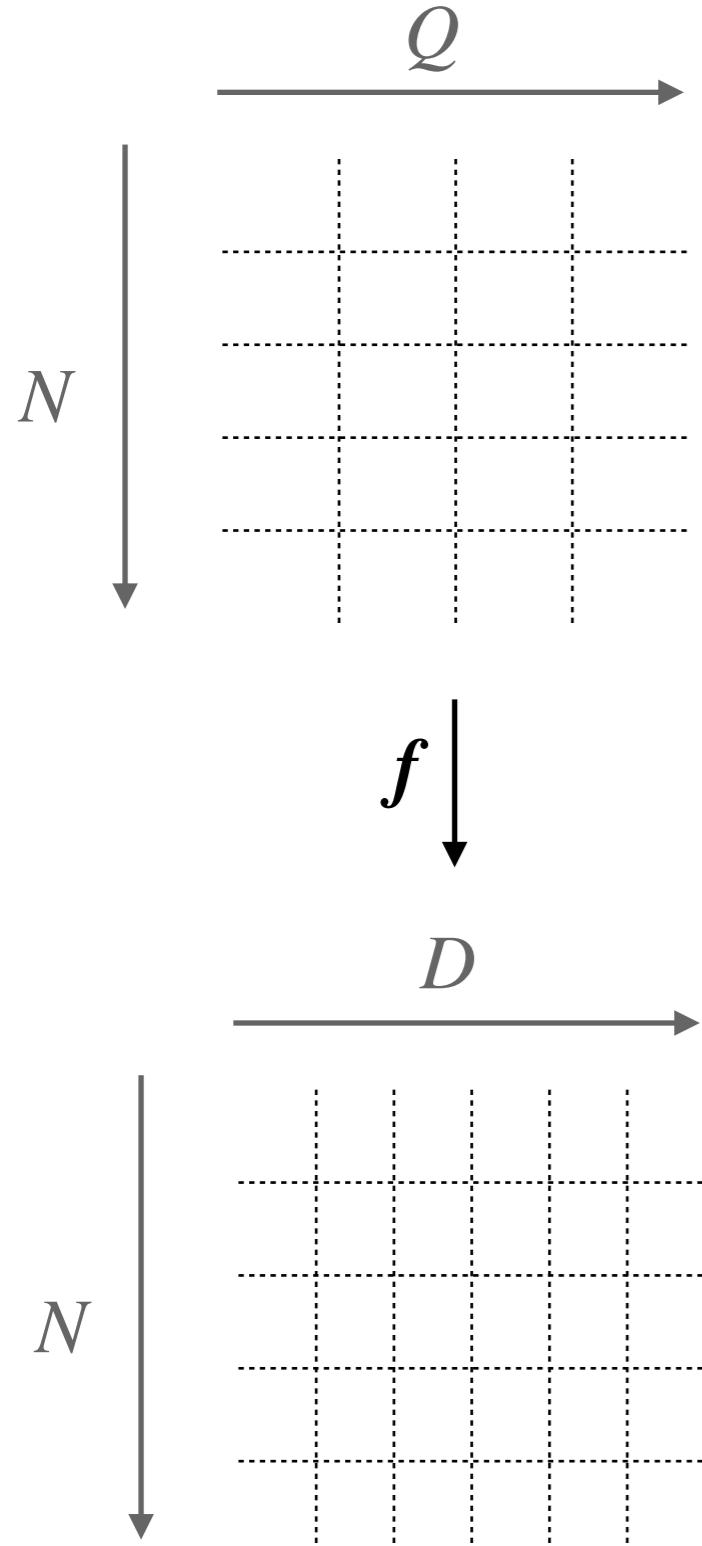




f

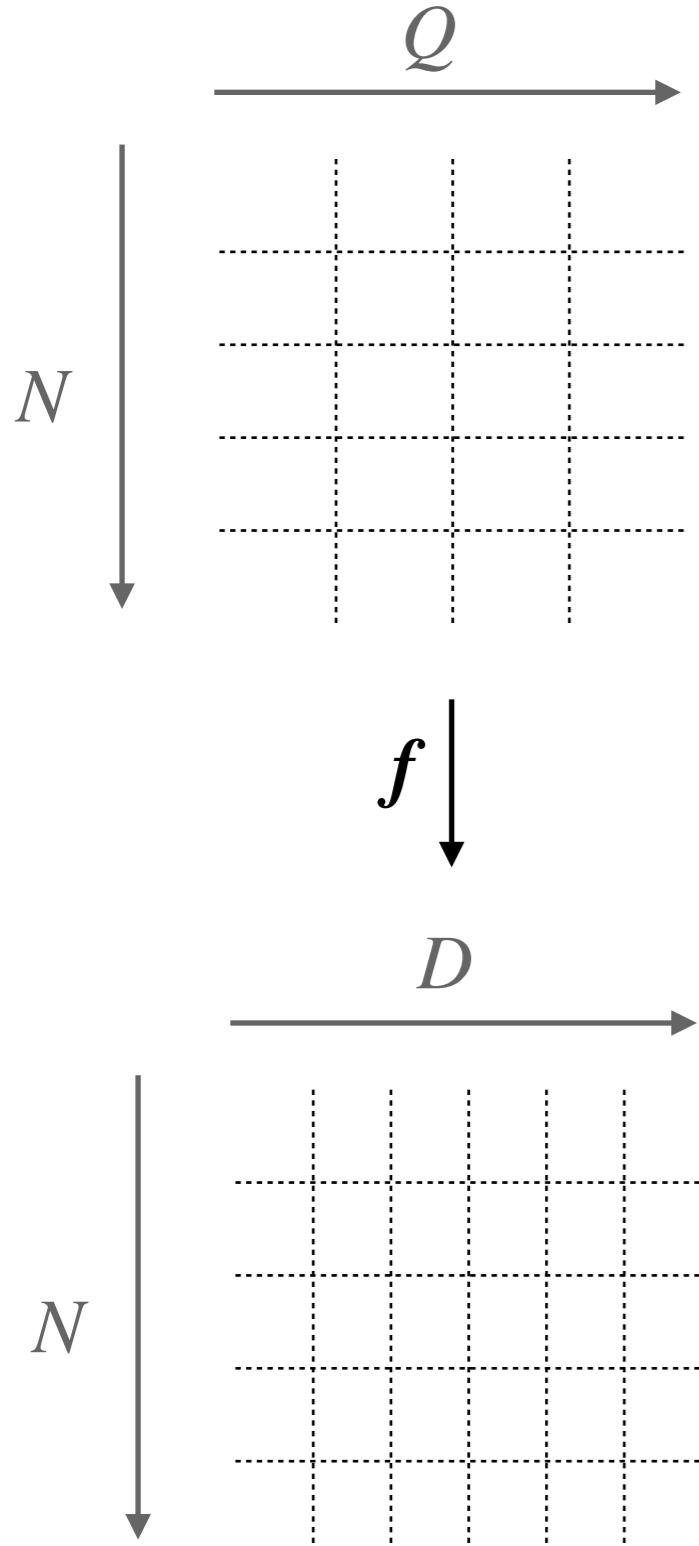


math
physics
language



$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix}$$

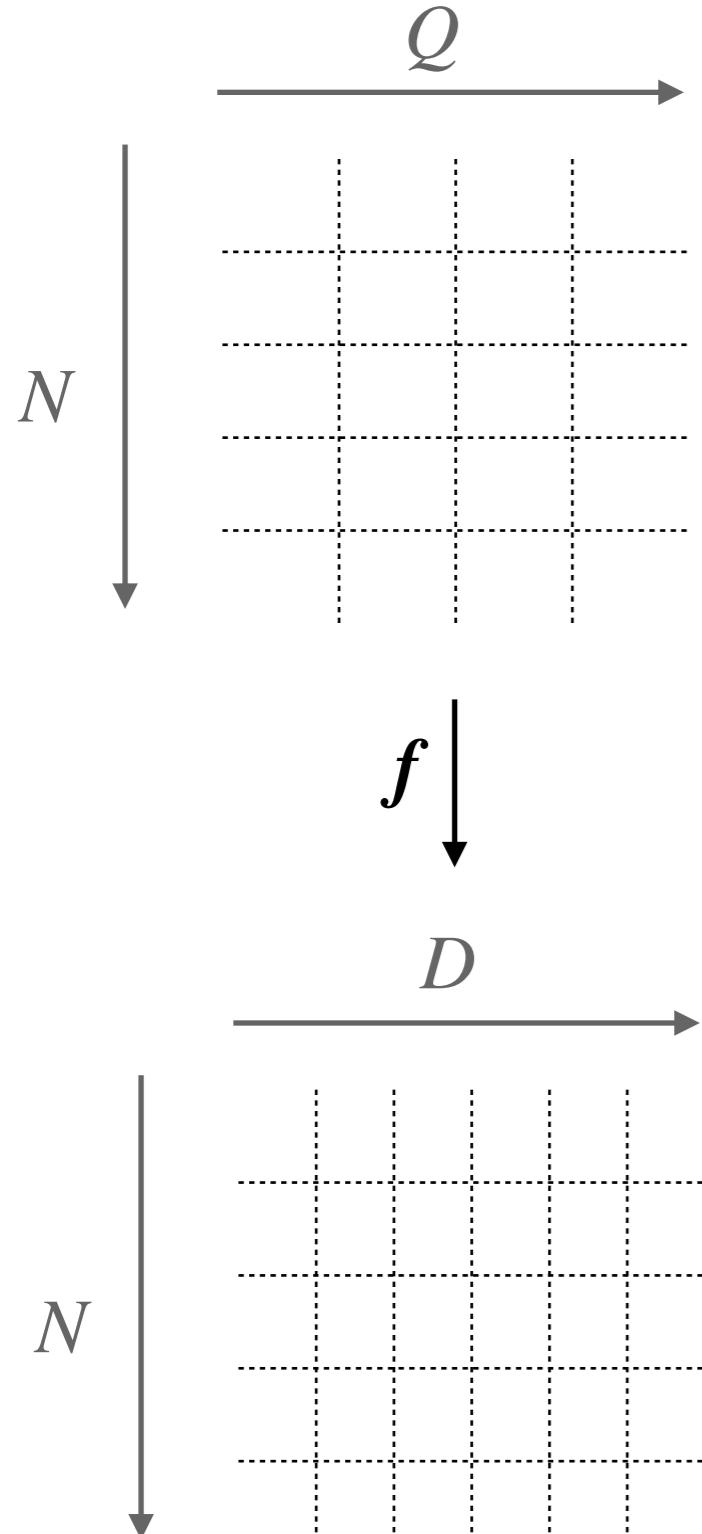
$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix} = f(x_1, x_2)$$



$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix}$$

$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix} = f(x_1, x_2)$$

$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix} = \begin{pmatrix} \uparrow & \downarrow \\ \uparrow & \downarrow \\ \downarrow & \uparrow \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$



$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix}$$

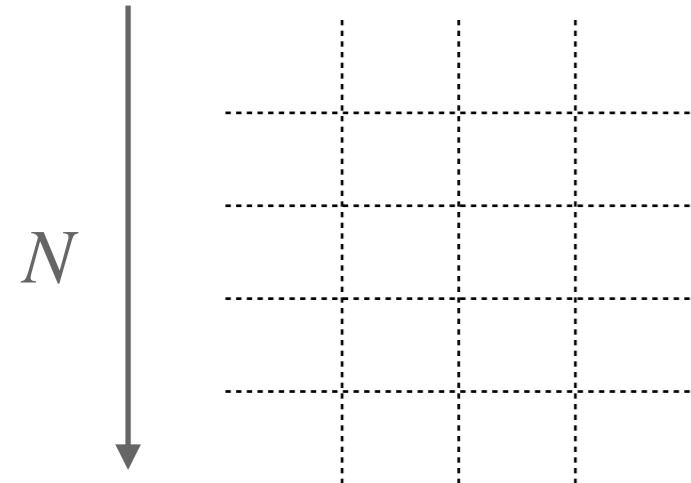
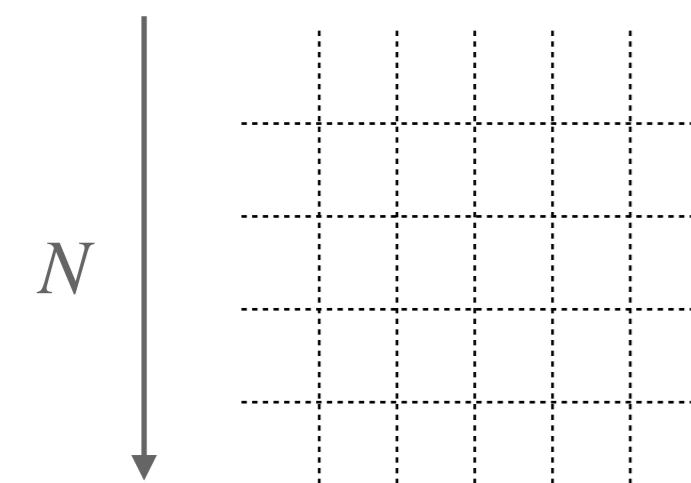
$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix} = f(x_1, x_2)$$

$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix} = \begin{pmatrix} \uparrow & \downarrow \\ \uparrow & \downarrow \\ \downarrow & \uparrow \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$\begin{pmatrix} \mathbf{math} \\ \mathbf{physics} \\ \mathbf{language} \end{pmatrix} = \begin{pmatrix} \uparrow & \downarrow \\ \uparrow & \downarrow \\ \downarrow & \uparrow \end{pmatrix} \begin{pmatrix} IQ \\ EQ \end{pmatrix}$$

Q

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

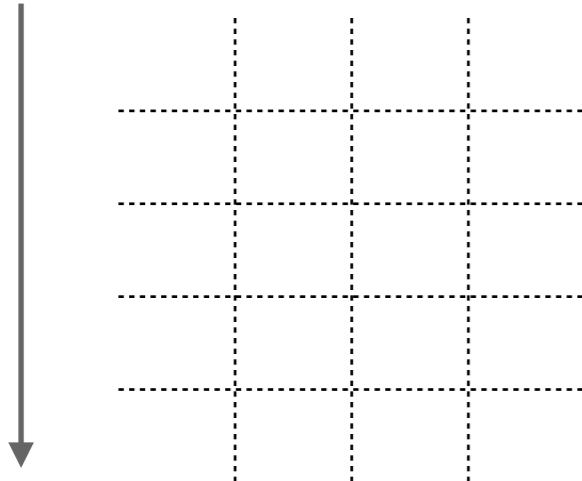
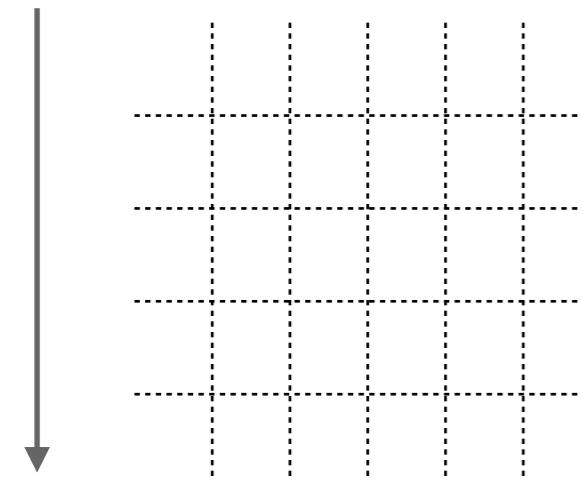
 f D 

Q

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

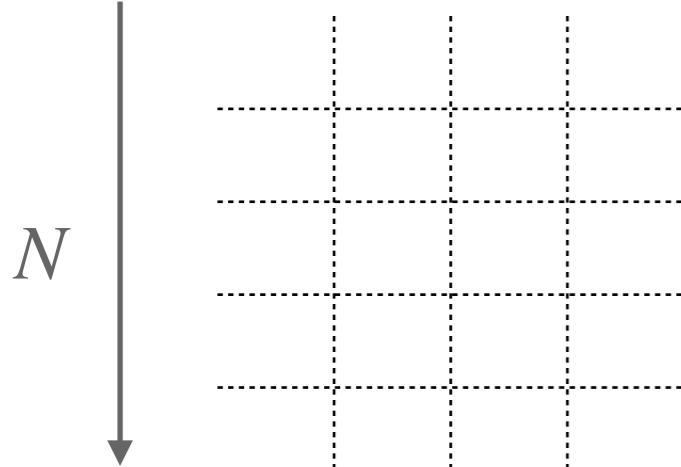
Generative point of view:

$$x \sim p_0(x)$$

 N  f D N 

Q

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



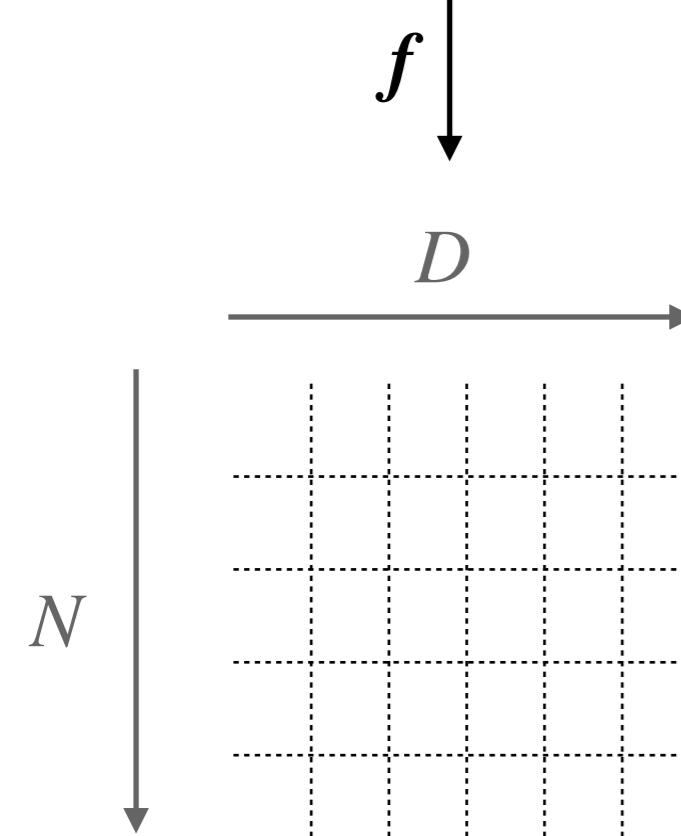
Generative point of view:

$$x \sim p_0(x)$$

$$y = f(x) + \epsilon$$

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X|Y)$ and $p(Y)$?

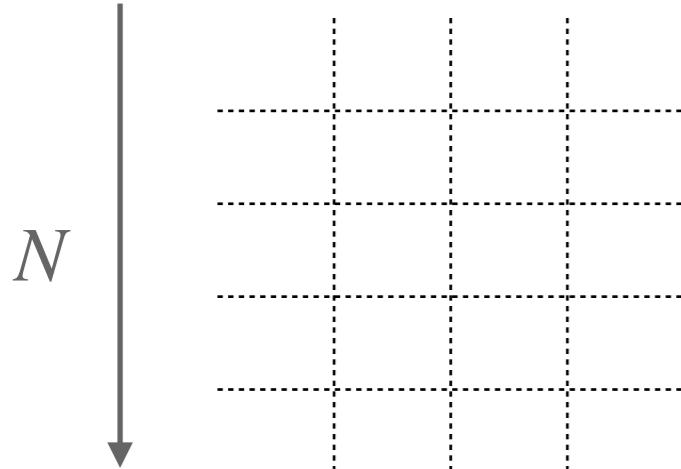
$$p(X|Y) = \frac{p(Y|X) p(X)}{p(Y)}$$



$$p(Y) = \int p(Y|X) p(X) dX$$

Q

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

**Generative point of view:**

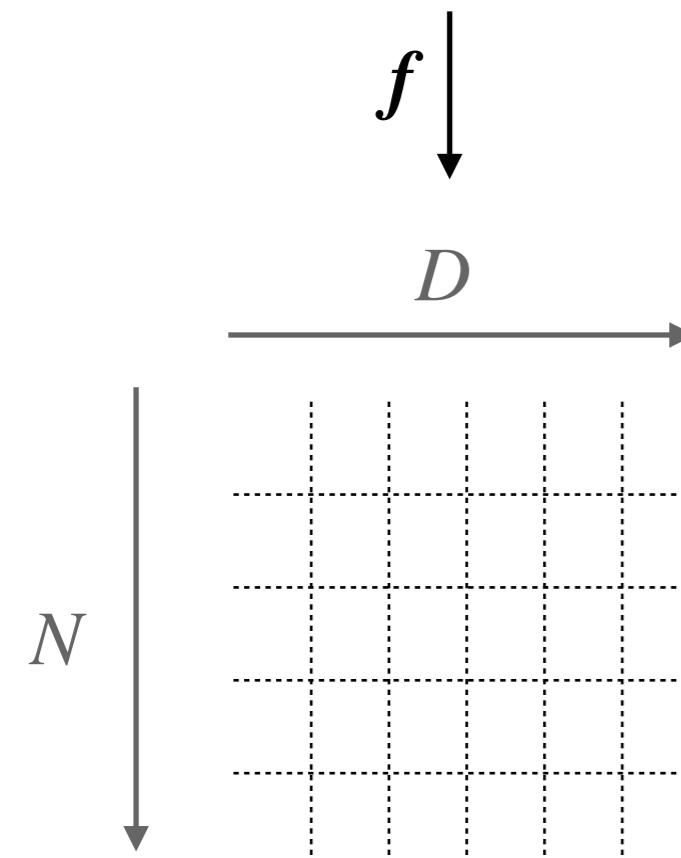
$$x \sim p_0(x)$$

$$y = f(x) + \epsilon$$

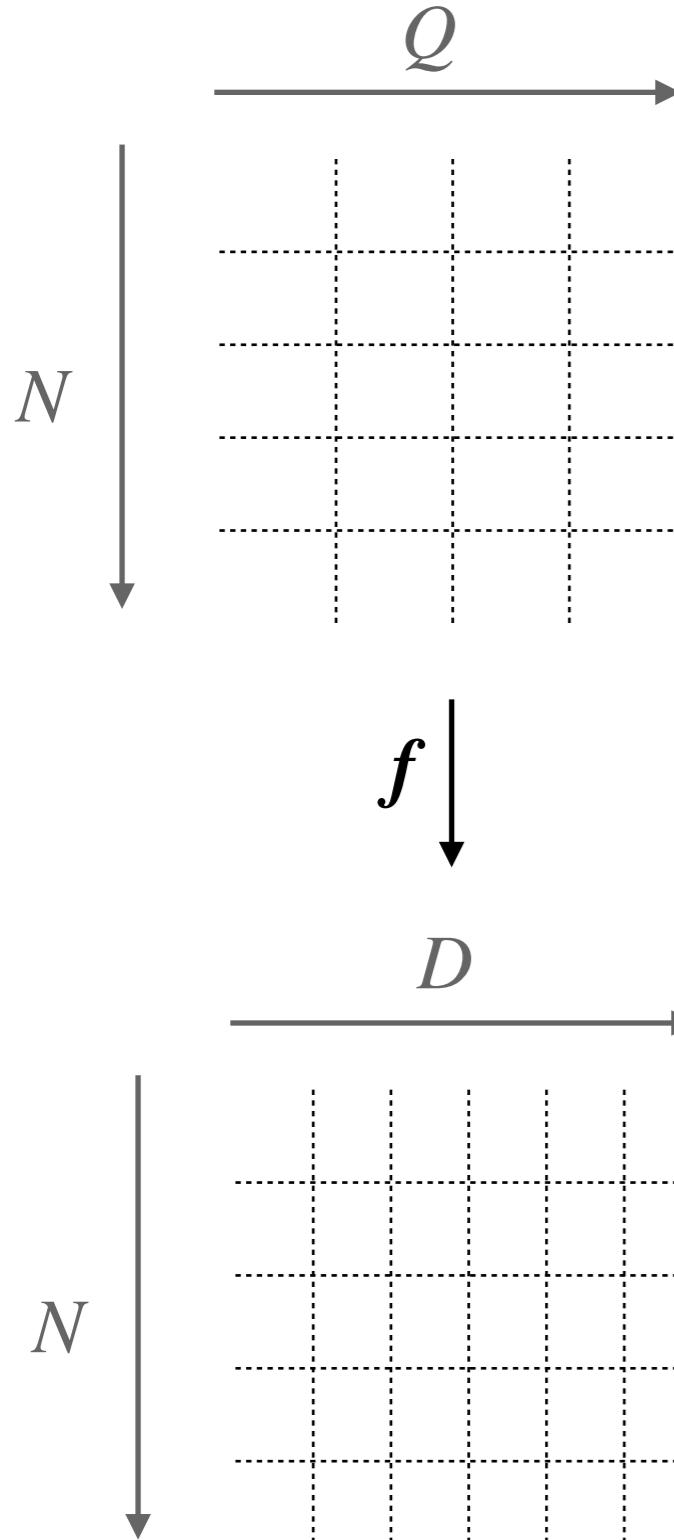
Given $Y = \{y_1, \dots, y_N\}$, what is $p(X|Y)$ and $p(Y)$?

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

$$p(Y) = \int p(Y|X) p(X) dX$$

**PPCA**

$$y = f(x) + \epsilon = Wx + \epsilon$$



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

Generative point of view:

$x \sim p_0(x)$

$y = f(x) + \epsilon$

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X|Y)$ and $p(Y)$?

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

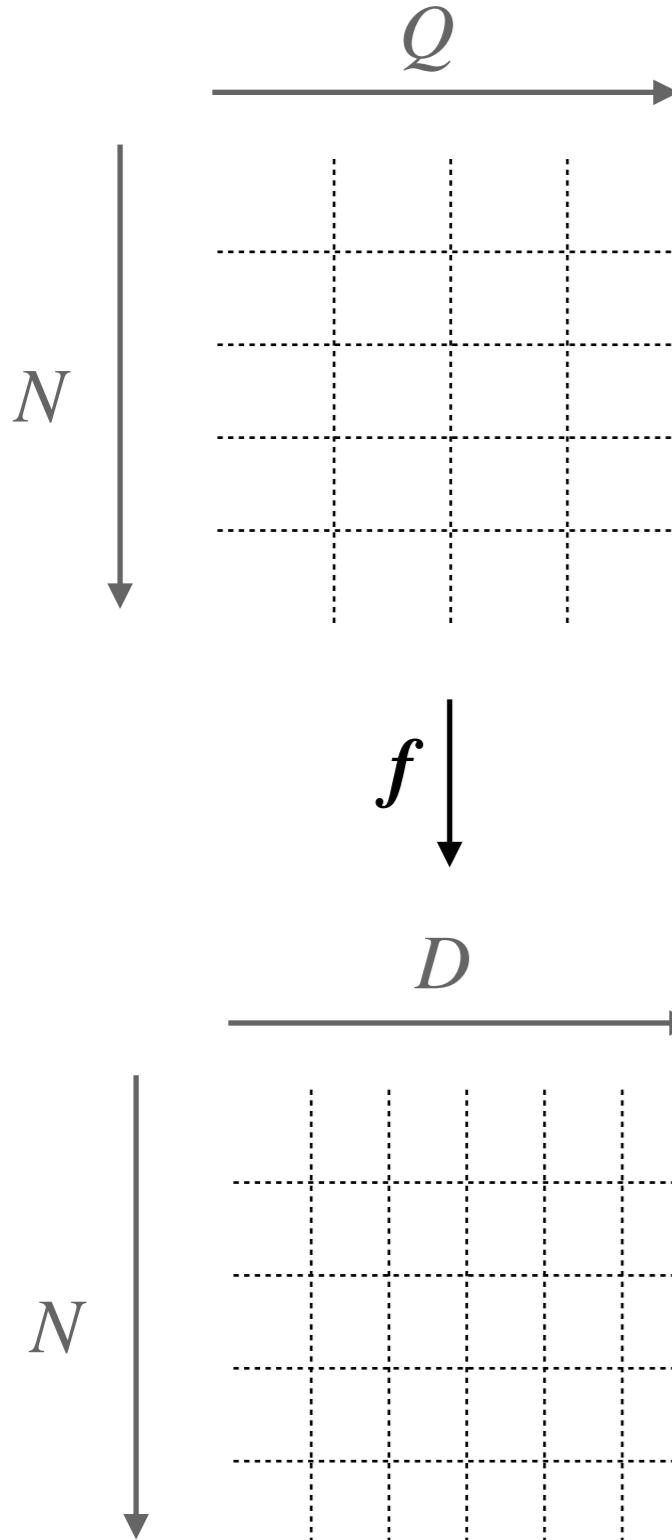
$$p(Y) = \int p(Y|X) p(X) dX$$

PPCA

$y = f(x) + \epsilon = Wx + \epsilon$

$p(Y|X, W) = \mathcal{N}(Y|XW^T, \sigma^2 I) = \prod_{n=1}^N \mathcal{N}(y_n | Wx_n, \sigma^2 I)$

$p(X) = \mathcal{N}(\mathbf{0}, I)$



Generative point of view:

$$x \sim p_0(x)$$

$$y = f(x) + \epsilon$$

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X|Y)$ and $p(Y)$?

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

$$p(Y) = \int p(Y|X) p(X) dX$$

PPCA

$$y = f(x) + \epsilon = Wx + \epsilon$$

$$p(Y|X, W) = \mathcal{N}(Y|XW^T, \sigma^2 I) = \prod_{n=1}^N \mathcal{N}(y_n | Wx_n, \sigma^2 I)$$

$$p(X) = \mathcal{N}(\mathbf{0}, I)$$

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

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

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

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X|Y)$ and $p(Y)$?

PPCA

$$y = f(x) + \epsilon = Wx + \epsilon$$

$$p(Y|X, W) = \mathcal{N}(Y|XW^T, \sigma^2 I)$$

$$= \prod_{n=1}^N \mathcal{N}(y_n | Wx_n, \sigma^2 I)$$

$$p(X) = \mathcal{N}(\mathbf{0}, I)$$

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

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

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X | Y)$ and $p(Y)$?

PPCA

$$y = f(x) + \epsilon = Wx + \epsilon$$

$$p(Y|X, W) = \mathcal{N}(Y|XW^T, \sigma^2 I)$$

$$= \prod_{n=1}^N \mathcal{N}(y_n | Wx_n, \sigma^2 I)$$

$$p(X) = \mathcal{N}(\mathbf{0}, I)$$

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

Variational Autoencoder

$$y = f_\theta(x) + \epsilon$$

$f_\theta(x)$ is NN

$$p(Y|X, \theta) = \prod_{n=1}^N \mathcal{N}(y_n | f_\theta(x_n), \sigma^2 I), \quad X \sim p(X)$$

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

$$p(Y) = \int p(Y|X) p(X) dX$$

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

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X | Y)$ and $p(Y)$?

PPCA

$$y = f(x) + \epsilon = Wx + \epsilon$$

$$p(Y|X, W) = \mathcal{N}(Y|XW^T, \sigma^2 I)$$

$$= \prod_{n=1}^N \mathcal{N}(y_n | Wx_n, \sigma^2 I)$$

$$p(X) = \mathcal{N}(\mathbf{0}, I)$$

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

Variational Autoencoder

$$y = f_\theta(x) + \epsilon$$

$f_\theta(x)$ is NN

$$p(Y|X, \theta) = \prod_{n=1}^N \mathcal{N}(y_n | f_\theta(x_n), \sigma^2 I), \quad X \sim p(X)$$

$$p(Y|\theta) = \int p(Y|X, \theta) p(X) dX$$

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

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X | Y)$ and $p(Y)$?

PPCA

$$y = f(x) + \epsilon = Wx + \epsilon$$

$$p(Y|X, W) = \mathcal{N}(Y|XW^T, \sigma^2 I)$$

$$= \prod_{n=1}^N \mathcal{N}(y_n | Wx_n, \sigma^2 I)$$

$$p(X) = \mathcal{N}(\mathbf{0}, I)$$

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

Variational Autoencoder

$$y = f_\theta(x) + \epsilon \quad f_\theta(x) \text{ is NN}$$

$$p(Y|X, \theta) = \prod_{n=1}^N \mathcal{N}(y_n | f_\theta(x_n), \sigma^2 I), \quad X \sim p(X)$$

$$\begin{aligned} p(Y|\theta) &= \int p(Y|X, \theta) p(X) dX \\ &= \int q(X) \frac{p(Y|X, \theta) p(X)}{q(X)} dX \end{aligned}$$

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

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

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X|Y)$ and $p(Y)$?

PPCA

$$y = f(x) + \epsilon = Wx + \epsilon$$

$$p(Y|X, W) = \mathcal{N}(Y|XW^T, \sigma^2 I)$$

$$= \prod_{n=1}^N \mathcal{N}(y_n | Wx_n, \sigma^2 I)$$

$$p(X) = \mathcal{N}(\mathbf{0}, I)$$

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

Variational Autoencoder

$$y = f_\theta(x) + \epsilon \quad f_\theta(x) \text{ is NN}$$

$$p(Y|X, \theta) = \prod_{n=1}^N \mathcal{N}(y_n | f_\theta(x_n), \sigma^2 I), \quad X \sim p(X)$$

$$\begin{aligned} p(Y|\theta) &= \int p(Y|X, \theta) p(X) dX \\ &= \int q(X) \frac{p(Y|X, \theta) p(X)}{q(X)} dX \end{aligned}$$

Meaningful $q(X)$ would be the posterior but is intractable \rightarrow Variational Approximation

$$q(X) = q_\phi(X|Y) = \mathcal{N}(x_n | f_\phi(y_n), \sigma^2 I)$$

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

Given $Y = \{y_1, \dots, y_N\}$, what is $p(X | Y)$ and $p(Y)$?

PPCA

$$y = f(x) + \epsilon = Wx + \epsilon$$

$$p(Y|X, W) = \mathcal{N}(Y|XW^T, \sigma^2 I)$$

$$= \prod_{n=1}^N \mathcal{N}(y_n | Wx_n, \sigma^2 I)$$

$$p(X) = \mathcal{N}(\mathbf{0}, I)$$

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

Variational Autoencoder

$$y = f_\theta(x) + \epsilon$$

$f_\theta(x)$ is NN

$$p(Y|X, \theta) = \prod_{n=1}^N \mathcal{N}(y_n | f_\theta(x_n), \sigma^2 I), \quad X \sim p(X)$$

$$p(Y|\theta) = \int p(Y|X, \theta) p(X) dX$$

$$= \int q(X) \frac{p(Y|X, \theta) p(X)}{q(X)} dX$$

Meaningful $q(X)$ would be the posterior but is intractable \rightarrow Variational Approximation

$$q(X) = q_\phi(X|Y) = \mathcal{N}(x_n | f_\phi(y_n), \sigma^2 I)$$

$$\ln p(Y|\theta) - \mathbf{KL} \left[q_\phi(X|Y) \| p(X|Y, \theta) \right] =$$

$$\mathbb{E}_{q_\phi(X|Y)} [\ln p(Y|X, \theta)] - \mathbf{KL} \left[q_\phi(X|Y) \| p(X) \right]$$

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

- 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}

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

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

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

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

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

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}

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

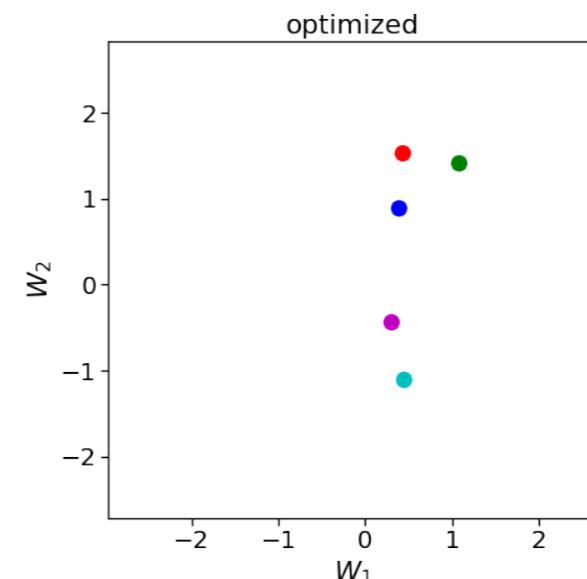
$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

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

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

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

- 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}

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

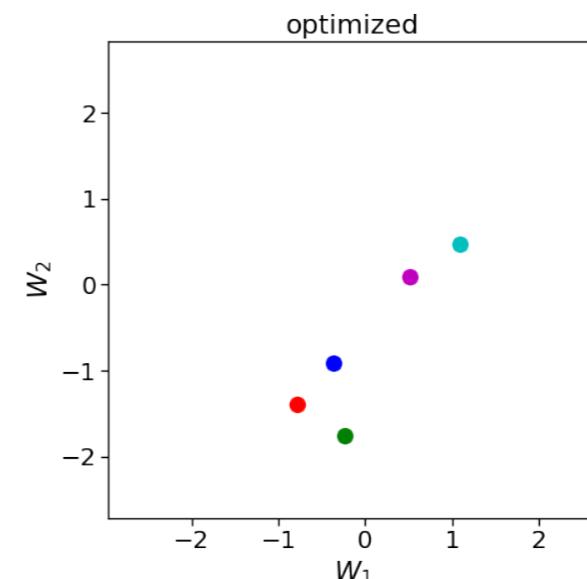
$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

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

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

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

- 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}

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

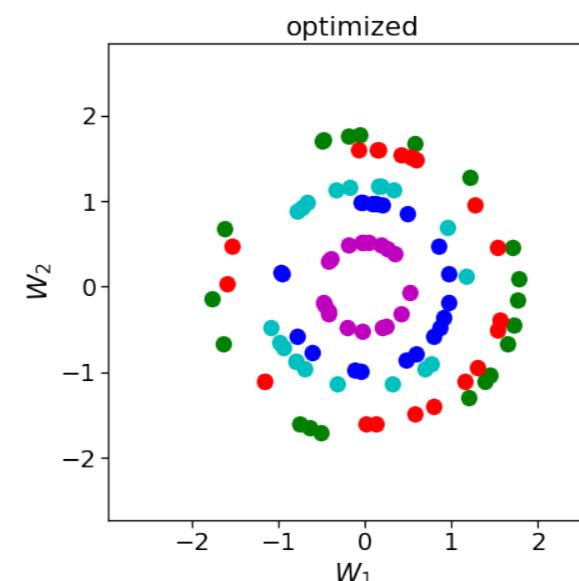
$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

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

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

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

- 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}

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

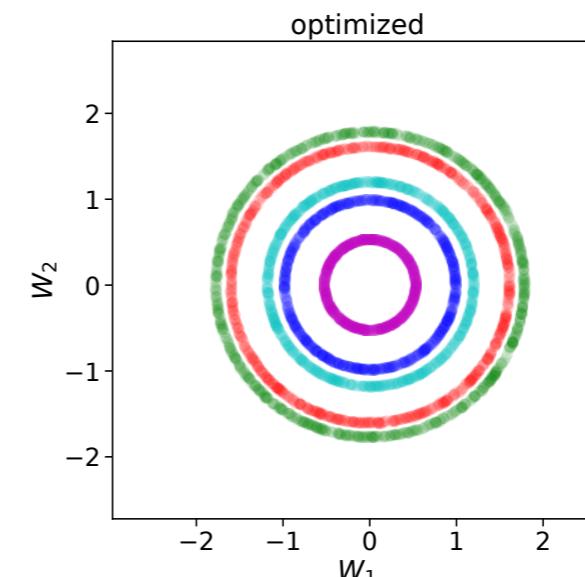
$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \epsilon$$

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

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

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

- Rotation invariant likelihood



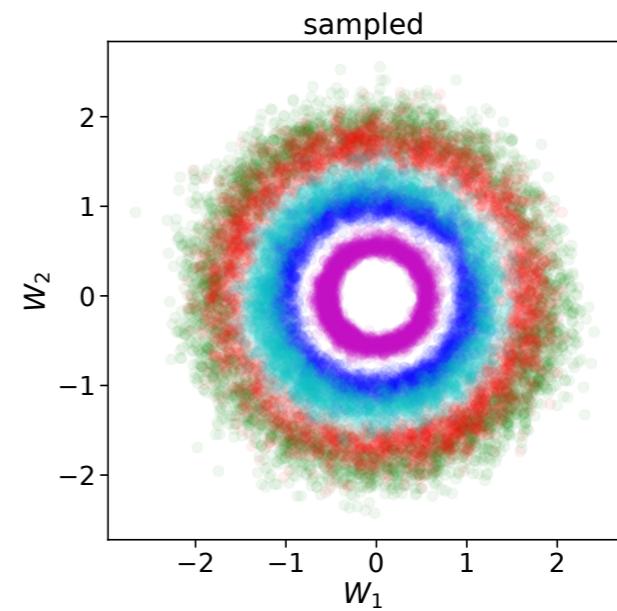
Bayesian approach to PPCA

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

- 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 \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
$$p(\mathbf{U}, \boldsymbol{\Sigma} | \mathbf{Y})$$

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**

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

Theory

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

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

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

→ **\boldsymbol{U} is uniformly distributed on the Stiefel manifold**

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

Theory

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

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

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 $\boldsymbol{\Sigma}$ is distributed as (James & Lee (2014))

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

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

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 $n = D : 1$

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

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

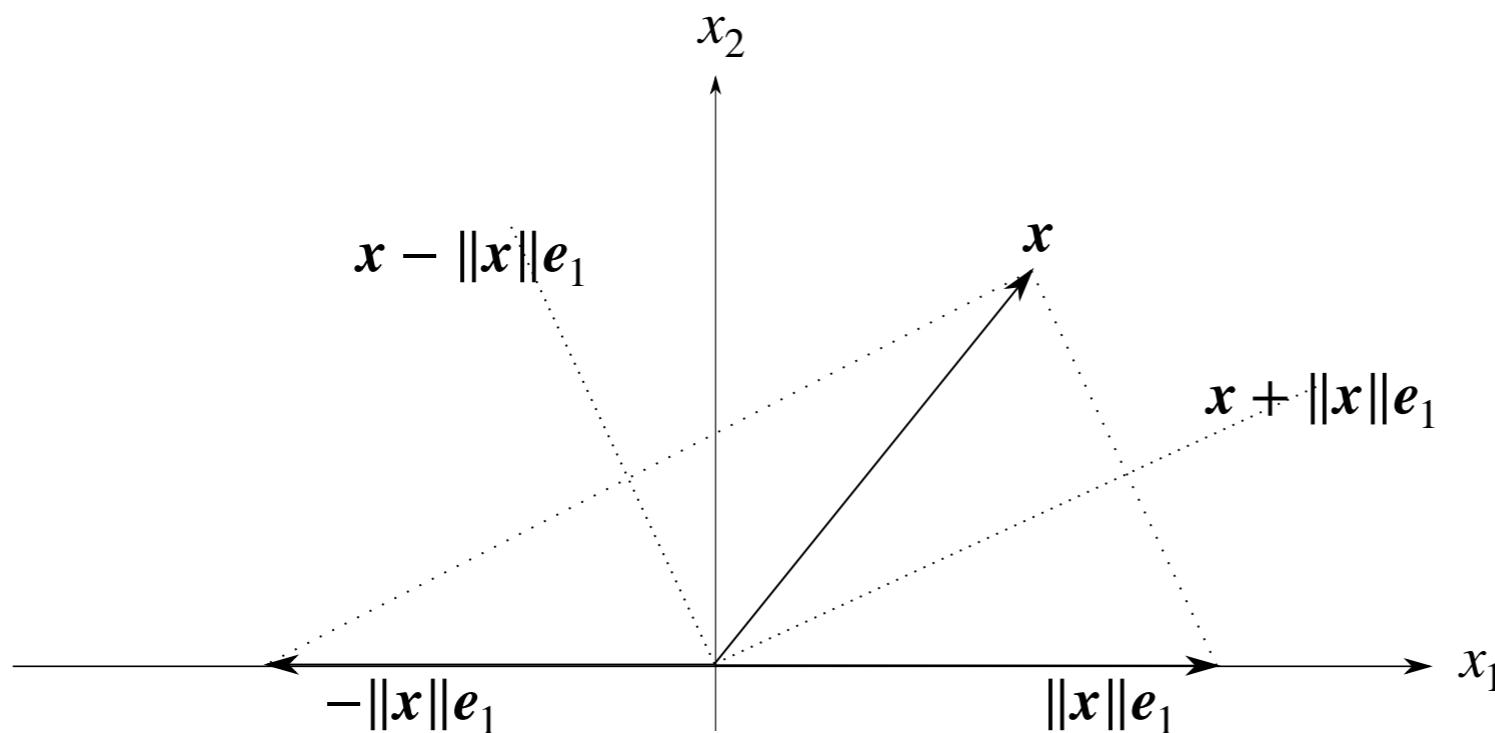
$$\tilde{H}_n(v_n) = -\text{sgn}(v_{n1})(I - 2u_n u_n^T)$$

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

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

Householder Transformations

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



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

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

Householder Transformations

Example for $(D, Q) = (2,2)$

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

for $n = D : 1$

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

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

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

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

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

Householder Transformations

Example for $(D, Q) = (2,2)$

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

Construction of \mathbf{H}_1

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

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

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

for $n = D : 1$

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

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

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

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

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

Householder Transformations

for $n = D : 1$

Example for $(D, Q) = (2, 2)$

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

$$\begin{aligned}\boldsymbol{v}_n &\sim \text{uniform on } \mathbb{S}^{n-1} \\ \boldsymbol{u}_n &= \frac{\boldsymbol{v}_n + \operatorname{sgn}(v_{n1}) \|\boldsymbol{v}_n\| \boldsymbol{e}_1}{\|\boldsymbol{v}_n + \operatorname{sgn}(v_{n1}) \|\boldsymbol{v}_n\| \boldsymbol{e}_1\|} \\ \tilde{\boldsymbol{H}}_n(\boldsymbol{v}_n) &= -\operatorname{sgn}(v_{n1})(\boldsymbol{I} - 2\boldsymbol{u}_n\boldsymbol{u}_n^T) \\ \boldsymbol{H}_n &= \begin{pmatrix} \boldsymbol{I} & \boldsymbol{0} \\ \boldsymbol{0} & \tilde{\boldsymbol{H}}_n \end{pmatrix}\end{aligned}$$

Construction of \boldsymbol{H}_1

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

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

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

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

Construction of \boldsymbol{H}_2

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

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

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

Householder Transformations

Example for $(D, Q) = (2,2)$

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

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

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

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

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

for $n = D : 1$

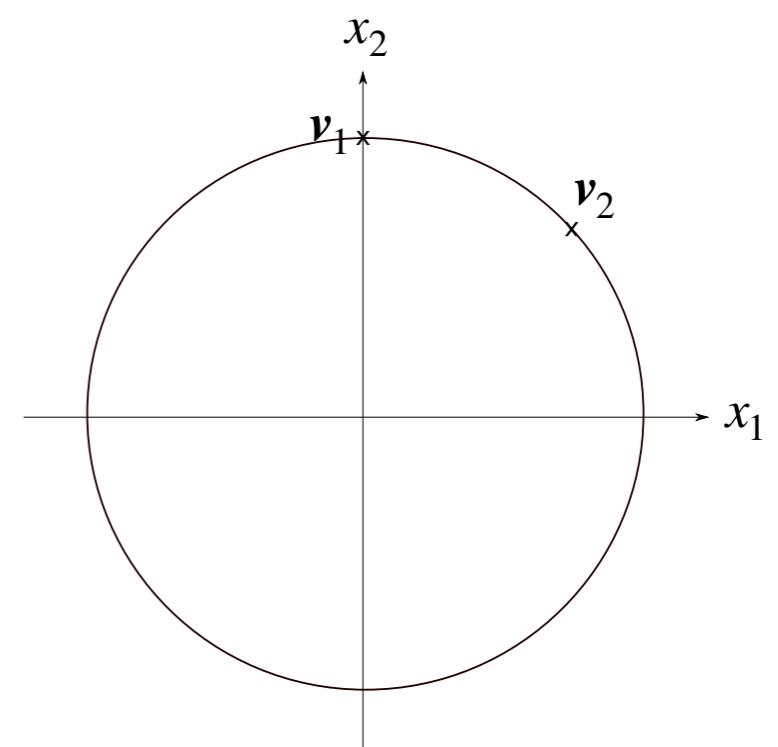
$$\boldsymbol{v}_n \sim \text{uniform on } \mathbb{S}^{n-1}$$

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

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

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

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



Householder Transformations

Example for $(D, Q) = (2,2)$

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

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

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

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

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

for $n = D : 1$

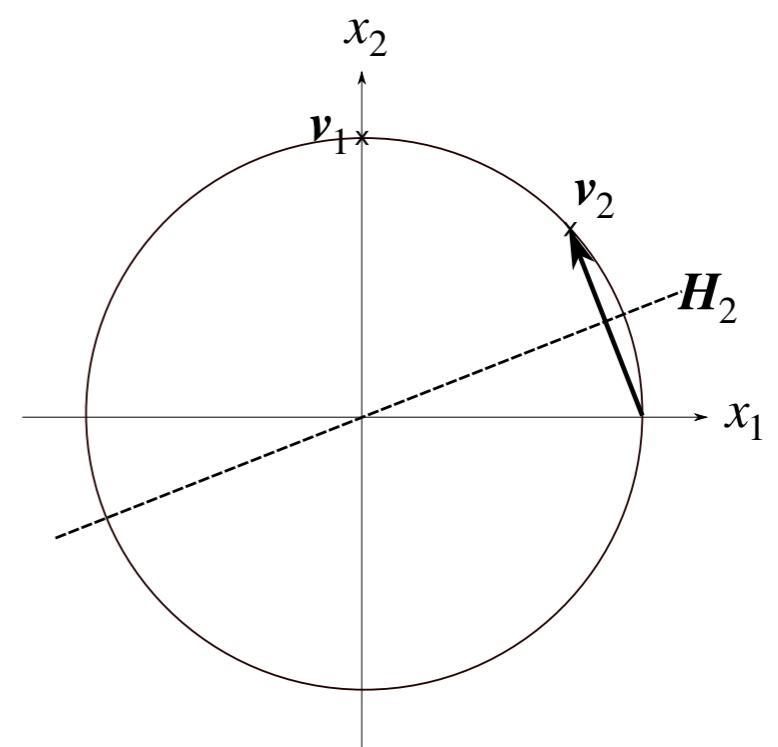
$\boldsymbol{v}_n \sim \text{uniform on } \mathbb{S}^{n-1}$

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

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

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

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



Householder Transformations

Example for $(D, Q) = (2,2)$

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

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

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

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

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

for $n = D : 1$

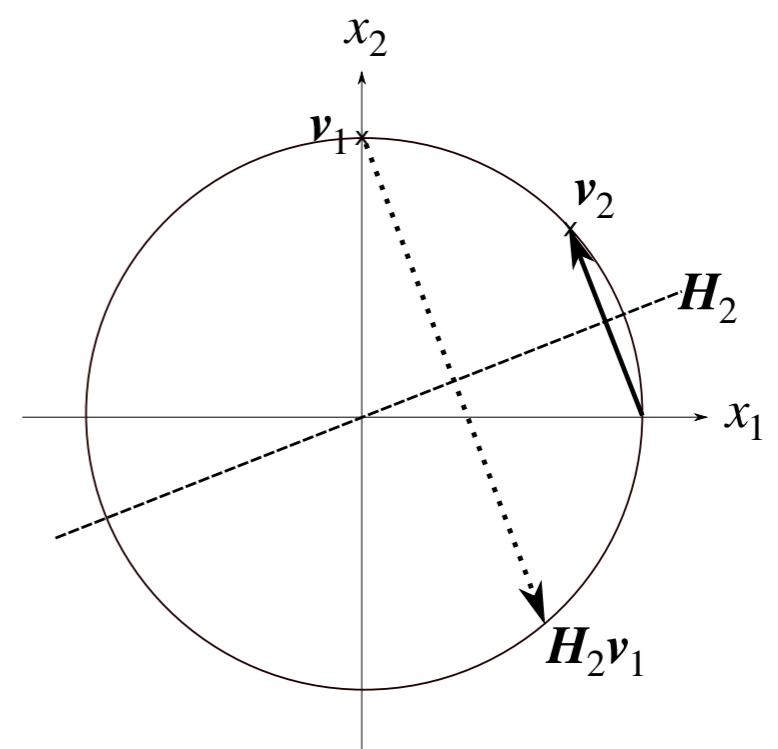
$\boldsymbol{v}_n \sim \text{uniform on } \mathbb{S}^{n-1}$

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

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

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

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



Householder Transformations

Example for $(D, Q) = (2,2)$

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

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

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

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

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

for $n = D : 1$

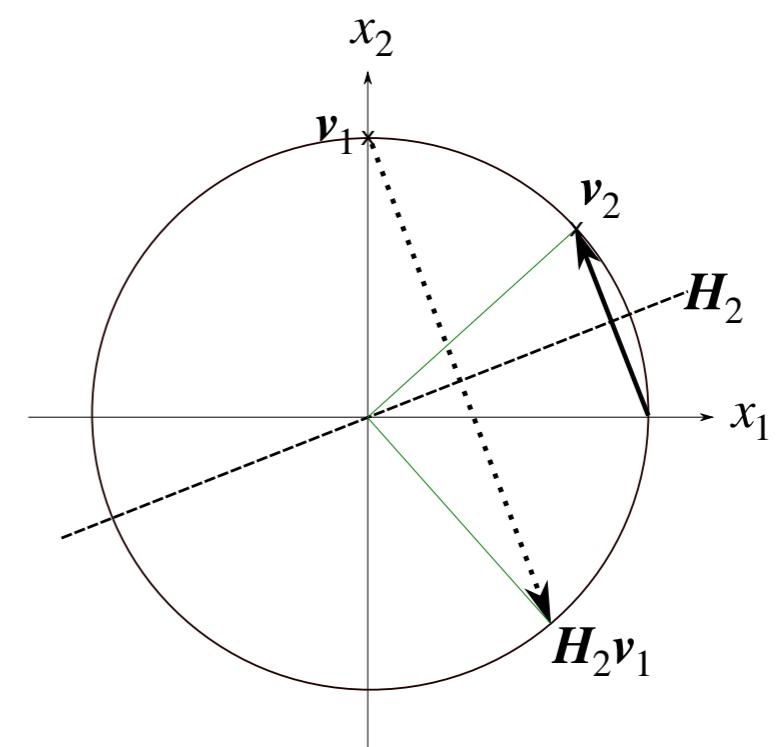
$\boldsymbol{v}_n \sim \text{uniform on } \mathbb{S}^{n-1}$

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

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

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

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



Householder Transformations

Example for $(D, Q) = (2,2)$

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

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

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

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

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

for $n = D : 1$

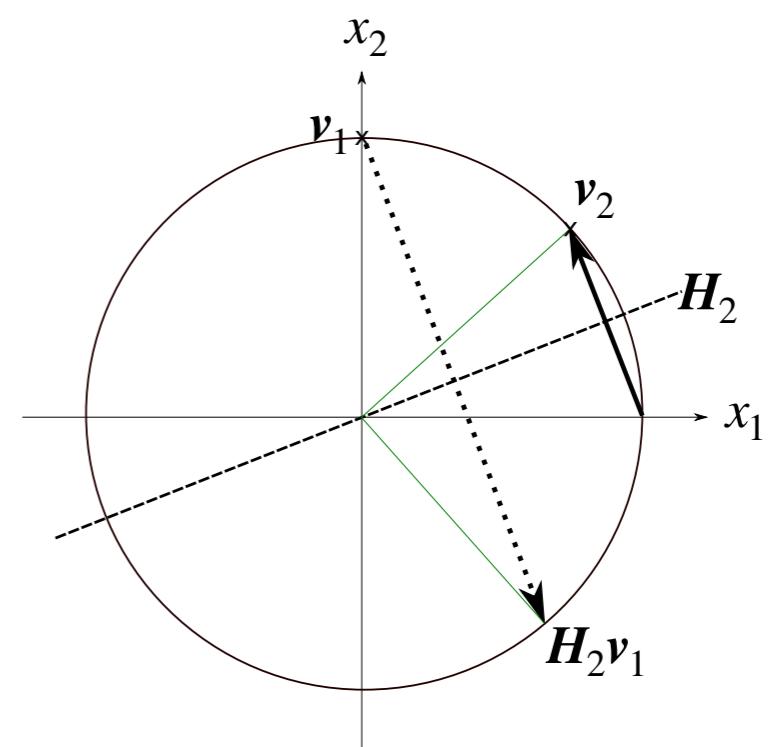
$\boldsymbol{v}_n \sim \text{uniform on } \mathbb{S}^{n-1}$

$$\mathbf{u}_n = \frac{\boldsymbol{v}_n + \text{sgn}(\boldsymbol{v}_{n1}) \|\boldsymbol{v}_n\| e_1}{\|\boldsymbol{v}_n + \text{sgn}(\boldsymbol{v}_{n1}) \|\boldsymbol{v}_n\| e_1\|}$$

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

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

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



Implementation

The full generative model for Bayesian PPCA:

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

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

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

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

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

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

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

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

Results

Synthetic Dataset

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

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

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

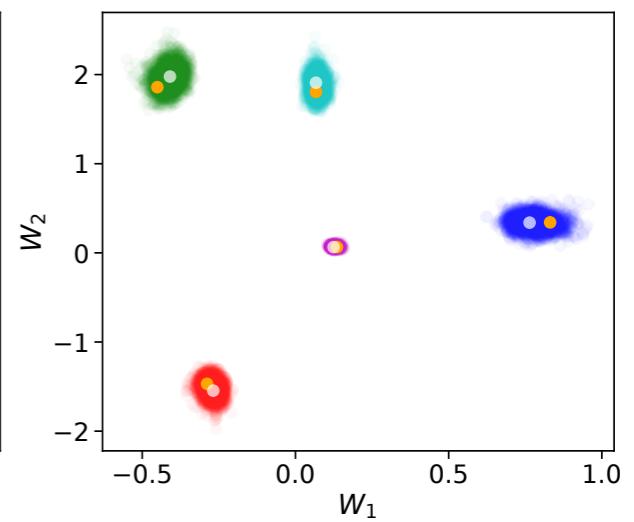
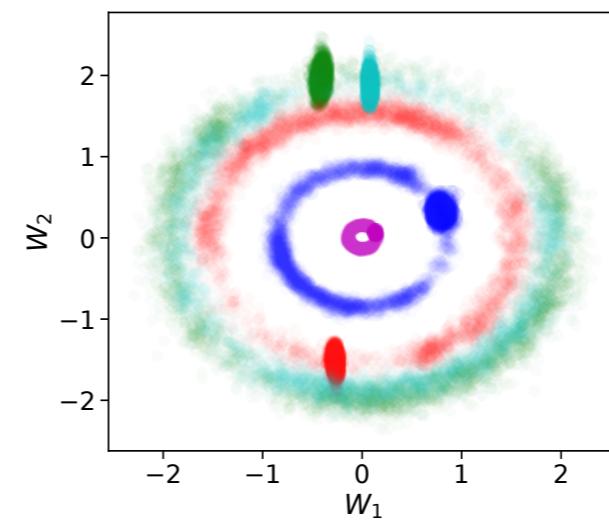
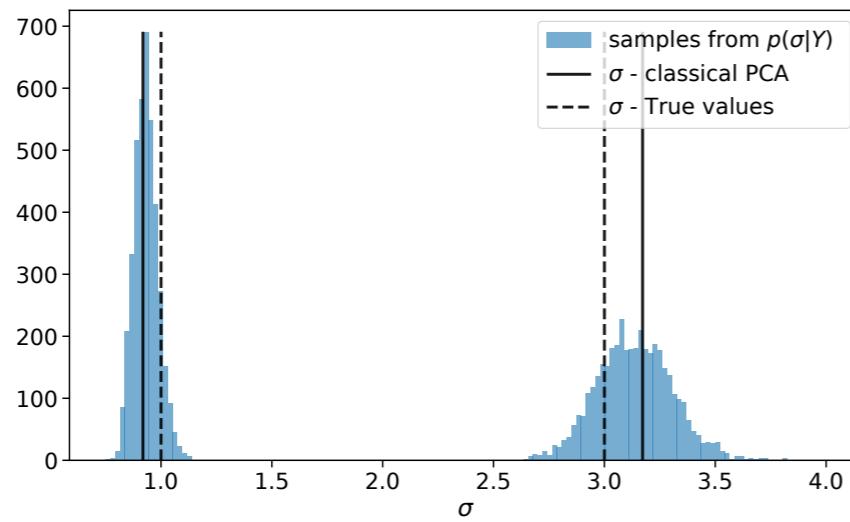
$$\epsilon \sim \mathcal{N}(0, 0.01) \in \mathbb{R}^{N \times D}$$

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

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

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

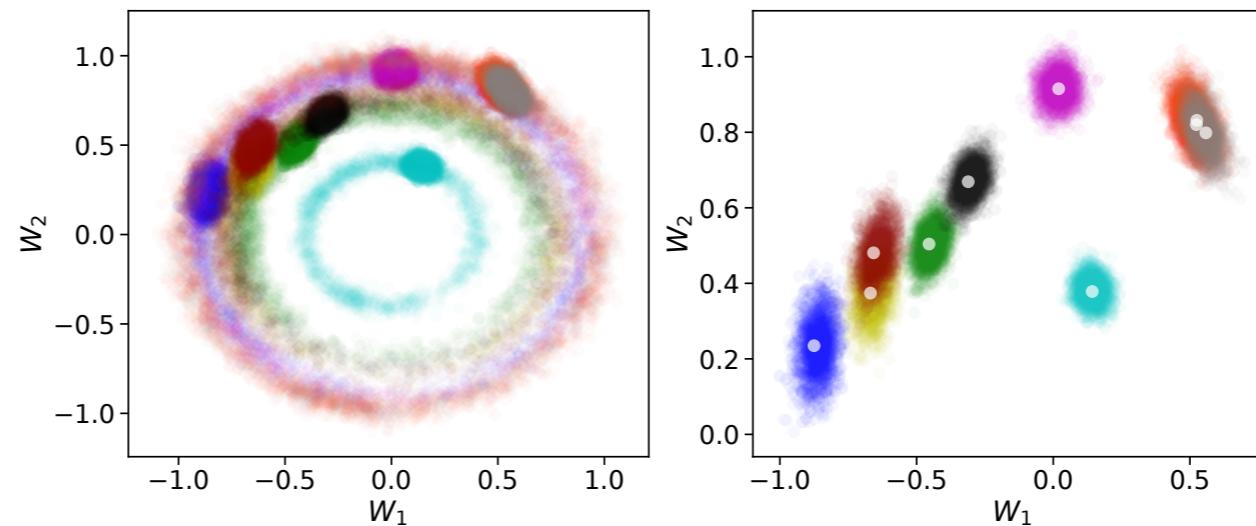
- Inference



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

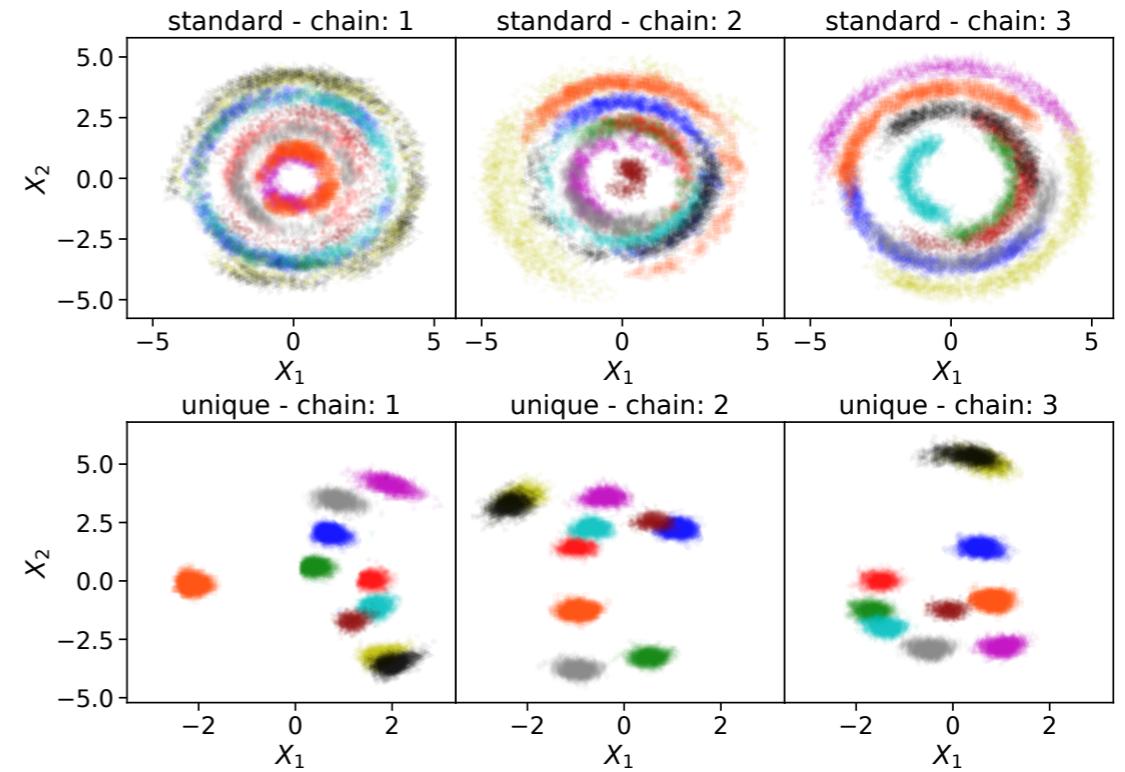
Extension to non-linear models

- 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 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**