# Probabilistic PCA

**Yingzhen Li**

Department of Computing
Imperial College London

@liyzhen2
yingzhen.li@imperial.ac.uk

Nov 25, 2022

## PCA: Recap

Motivation: real-world data $\mathcal{D} = \{x_n\}_{n=1}^N, x_n \in \mathbb{R}^{D \times 1}$ often lies in a lower-dimensional space

PCA's idea to "save memory":

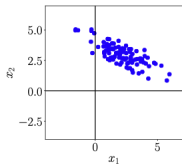- Project $x_n$ onto a lower-dim space $span(\{\mathbf{b}_1, ..., \mathbf{b}_M\})$ to get
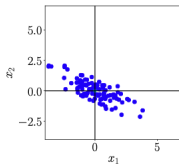
$$z_n := (z_{n1}, ..., z_{nM}), \quad z_{nm} = \mathbf{b}_m, \quad M < D,$$

  then store $z_n$ instead of $x_n$;

- When needed, get reconstruction $\tilde{x}_n = \sum_{m=1}^M z_{nm} \mathbf{b}_m$

- To get orthonormal basis $\{\mathbf{b}_1, ..., \mathbf{b}_M\}$: PCA
  - maximum variance view
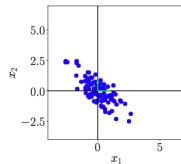  - minimum reconstruction error view

# PCA: Recap



(a) Original dataset.

(b) Step 1: Centering by subtracting the mean from each data point.

(c) Step 2: Dividing by the standard deviation to make the data unit free. Data has variance 1 along each axis.

(d) Step 3: Compute eigenvalues and eigenvectors (arrows) of the data covariance matrix (ellipse).

(e) Step 4: Project data onto the principal subspace.

(f) Undo the standardization and move projected data back into the original data space from (a).

Fig from the MML book.

# PCA: Recap

An issue with PCA in test time:

- Given an $x$, we can find the low-dim projection $z$ of it using trained PCA
- However, PCA alone cannot generate new $x$ (unless we do something further)

# Generative models

To name a few dimensionality reduction methods:

# Generative models



Sprouts in the shape of text 'Imagen' coming out of a fairytale book.

A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.

A high contrast portrait of a very happy fuzzy panda dressed as a chef in a high end kitchen making dough. There is a painting of flowers on the wall behind him.

Teddy bears swimming at the Olympics 400m Butterfly event.

A cute corgi lives in a house made out of sushi.

A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.

Google's Imagen text-to-image generation model

# Latent Variable Models

Data distribution: $\mathcal{D} = \{x_n\}_{n=1}^N, x_n \sim \pi(x)$
Make a generative model that generates $x$ as follows:

$$z \sim p_\theta(z), \quad x \sim p_\theta(x|z)$$

- $z$: latent variable

- $x$: data

- $\theta$: model parameter to be fitted

- if $p_\theta(x) \approx \pi(x)$, then the model can generate realistic data

# Probabilistic PCA

Data distribution: $\mathcal{D} = \{x_n\}_{n=1}^N, x_n \sim \pi(x)$
Probabilistic PCA: make a latent variable model as follows:

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

$$p_{\theta}(x|z) = \mathcal{N}(x; \mathbf{W}z + \mu, \sigma^2\mathbf{I})$$

Sampling from this generative model:

$$z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad x = \mathbf{W}z + \mu + \sigma\epsilon, \ \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

‣ model parameter: $\theta = \{\mathbf{W}, \mu\}$

# Probabilistic PCA

Data distribution: $\mathcal{D} = \{x_n\}_{n=1}^N$, $x_n \sim \pi(x)$
Probabilistic PCA: make a latent variable model as follows:

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

$$p_\theta(x|z) = \mathcal{N}(x; \mathbf{W}z + \mu, \sigma^2 \mathbf{I})$$

Marginal distribution:

$$p_\theta(x) = \int p_\theta(x|z)p(z)dz$$

# Probabilistic PCA

Data distribution: $\mathcal{D} = \{x_n\}_{n=1}^N, x_n \sim \pi(x)$
Probabilistic PCA: make a latent variable model as follows:

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

$$p_\theta(x|z) = \mathcal{N}(x; \mathbf{W}z + \mu, \sigma^2\mathbf{I})$$

Marginal distribution:

$$p_\theta(x) = \int p_\theta(x|z)p(z)dz$$
$$= \mathcal{N}(x; \mu, \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}).$$

# Probabilistic PCA

Data distribution: $\mathcal{D} = \{x_n\}_{n=1}^{N}, x_n \sim \pi(x)$

Probabilistic PCA: make a latent variable model as follows:

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

$$p_{\theta}(x|z) = \mathcal{N}(x; \mathbf{W}z + \mu, \sigma^2 \mathbf{I})$$

Fitting $\theta$ with Maximum Likelihood Estimation (MLE):

$$\max_{\theta} \mathcal{L}(\theta), \quad \mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^{N} \log p_{\theta}(x_n)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \log \mathcal{N}(x_n; \mu, \mathbf{W}\mathbf{W}^{\top} + \sigma^2 \mathbf{I})$$

# Probabilistic PCA

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}), \quad \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^{N} \log \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}), \quad \boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\mu}\}$$

Derivative of $\mathcal{L}$ w.r.t. $\boldsymbol{\mu}$: denoting $\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}$

$$\frac{\partial}{\partial \boldsymbol{\mu}} \log \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}) = \frac{\partial}{\partial \boldsymbol{\mu}} \left( -\frac{1}{2}(\boldsymbol{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1}(\boldsymbol{x}_n - \boldsymbol{\mu}) \right)$$

# Probabilistic PCA

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}), \quad \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^{N} \log \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}), \quad \boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\mu}\}$$

Derivative of $\mathcal{L}$ w.r.t. $\boldsymbol{\mu}$: denoting $\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}$

$$\frac{\partial}{\partial \boldsymbol{\mu}} \log \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}) = \frac{\partial}{\partial \boldsymbol{\mu}} \left( -\frac{1}{2} (\boldsymbol{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\boldsymbol{x}_n - \boldsymbol{\mu}) \right)$$
$$= (\boldsymbol{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1}$$

$$\Rightarrow \quad \frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^{N} (\boldsymbol{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1}$$

Setting $\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}} = \mathbf{0} \quad \Rightarrow \quad \boldsymbol{\mu}^* = \frac{1}{N} \sum_{n=1}^{N} \boldsymbol{x}_n$

# Probabilistic PCA

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}), \quad \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^{N} \log \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}), \quad \boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\mu}, \sigma\}$$

Derivative of $\mathcal{L}$ w.r.t. $\mathbf{W}$: denoting $\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}$

$$\frac{\partial}{\partial \mathbf{W}} \log \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})$$
$$= \frac{\partial}{\partial \mathbf{W}} \left( -\frac{1}{2}(\boldsymbol{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\boldsymbol{x}_n - \boldsymbol{\mu}) - \frac{1}{2} \log |\mathbf{C}| \right)$$

- $\mathbf{C}$ depends on $\mathbf{W}$, so "chain rule" applies
- However, so far we've only learned about chain rule applied to scalars and vectors.

# Probabilistic PCA

Applying chain rule: let $\mathcal{L}_n := \log \mathcal{N}(x_n; \mu, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})$

‣ Chain rule for individual elements $W_{kl}$ of $\mathbf{W}$:

$$\frac{\partial \mathcal{L}_n}{\partial W_{kl}} = \sum_{i,j} \frac{\partial \mathcal{L}_n}{\partial C_{ij}} \frac{\partial C_{ij}}{\partial W_{kl}}, \quad \mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}$$

‣ Calculate $\frac{\partial C_{ij}}{\partial W_{kl}}$:

# Probabilistic PCA

Applying chain rule: let $\mathcal{L}_n := \log \mathcal{N}(x_n; \mu, \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I})$

‣ Chain rule for individual elements $W_{kl}$ of $\mathbf{W}$:

$$\frac{\partial \mathcal{L}_n}{\partial W_{kl}} = \sum_{i,j} \frac{\partial \mathcal{L}_n}{\partial C_{ij}} \frac{\partial C_{ij}}{\partial W_{kl}}, \quad \mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}$$

‣ Calculate $\frac{\partial C_{ij}}{\partial W_{kl}}$: notice $C_{ij} = \sum_l W_{il} W_{jl} + \sigma^2 \delta(i = j)$

$$\frac{\partial C_{ij}}{\partial W_{kl}} = \begin{cases} 0, & k \notin \{i, j\} \\ W_{jl}, & k = i \neq j \\ W_{il}, & k = j \neq i \\ 2W_{il}, & k = i = j \end{cases}$$

# Probabilistic PCA

Applying chain rule: let $\mathcal{L}_n := \log \mathcal{N}(x_n; \mu, WW^\top + \sigma^2 I)$

‣ Chain rule for individual elements $W_{kl}$ of $W$:

$$\frac{\partial \mathcal{L}_n}{\partial W_{kl}} = \sum_{i,j} \frac{\partial \mathcal{L}_n}{\partial C_{ij}} \frac{\partial C_{ij}}{\partial W_{kl}}, \quad C = WW^\top + \sigma^2 I$$

‣ Calculate $\frac{\partial C_{ij}}{\partial W_{kl}}$: notice $C_{ij} = \sum_l W_{il} W_{jl} + \sigma^2 \delta(i = j)$

$$\frac{\partial C_{ij}}{\partial W_{kl}} = \begin{cases} 0, & k \notin \{i, j\} \\ W_{jl}, & k = i \neq j \\ W_{il}, & k = j \neq i \\ 2W_{il}, & k = i = j \end{cases}$$

‣ This means for fixed $k, l$:

$$\frac{\partial \mathcal{L}_n}{\partial W_{kl}} = \sum_j \frac{\partial \mathcal{L}_n}{\partial C_{kj}} W_{jl} + \sum_i \frac{\partial \mathcal{L}_n}{\partial C_{ik}} W_{il}$$

# Probabilistic PCA

Applying chain rule: let $\mathcal{L}_n := \log \mathcal{N}(x_n; \mu, WW^\top + \sigma^2 I)$

- Chain rule for individual elements $W_{kl}$ of $W$:

$$\frac{\partial \mathcal{L}_n}{\partial W_{kl}} = \sum_j \frac{\partial \mathcal{L}_n}{\partial C_{kj}} W_{jl} + \sum_i \frac{\partial \mathcal{L}_n}{\partial C_{ik}} W_{il}$$

- Writing the derivatives of $\mathcal{L}_n$ in matrix forms:

$$\frac{\partial \mathcal{L}_n}{\partial C} = \begin{bmatrix} \frac{\partial \mathcal{L}_n}{\partial C_{11}} & \frac{\partial \mathcal{L}_n}{\partial C_{21}} & \cdots \\ \vdots & \ddots & \\ \frac{\partial \mathcal{L}_n}{\partial C_{1D}} & & \frac{\partial \mathcal{L}_n}{\partial C_{DD}} \end{bmatrix}, \quad \frac{\partial \mathcal{L}_n}{\partial W} = \begin{bmatrix} \frac{\partial \mathcal{L}_n}{\partial W_{11}} & \frac{\partial \mathcal{L}_n}{\partial W_{21}} & \cdots \\ \vdots & \ddots & \\ \frac{\partial \mathcal{L}_n}{\partial W_{1M}} & & \frac{\partial \mathcal{L}_n}{\partial W_{DM}} \end{bmatrix}$$

$$\Rightarrow \quad \sum_j \frac{\partial \mathcal{L}_n}{\partial C_{kj}} W_{jl} = (W^\top)_{l\cdot} \left( \frac{\partial \mathcal{L}_n}{\partial C} \right)_{\cdot k}, \quad \sum_i \frac{\partial \mathcal{L}_n}{\partial C_{ik}} W_{il} = \left( \frac{\partial \mathcal{L}_n}{\partial C} \right)_{k\cdot} W_{\cdot l}$$

## Probabilistic PCA

Applying chain rule: let $\mathcal{L}_n := \log \mathcal{N}(x_n; \mu, WW^\top + \sigma^2 I)$

‣ Chain rule for individual elements $W_{kl}$ of $W$:

$$\frac{\partial \mathcal{L}_n}{\partial W_{kl}} = \sum_j \frac{\partial \mathcal{L}_n}{\partial C_{kj}} W_{jl} + \sum_i \frac{\partial \mathcal{L}_n}{\partial C_{ik}} W_{il}$$

‣ Writing the derivatives of $\mathcal{L}_n$ in matrix forms:

$$\frac{\partial \mathcal{L}_n}{\partial C} = \begin{bmatrix} \frac{\partial \mathcal{L}_n}{\partial C_{11}} & \frac{\partial \mathcal{L}_n}{\partial C_{21}} & \cdots \\ \vdots & \ddots & \\ \frac{\partial \mathcal{L}_n}{\partial C_{1D}} & & \frac{\partial \mathcal{L}_n}{\partial C_{DD}} \end{bmatrix}, \quad \frac{\partial \mathcal{L}_n}{\partial W} = \begin{bmatrix} \frac{\partial \mathcal{L}_n}{\partial W_{11}} & \frac{\partial \mathcal{L}_n}{\partial W_{21}} & \cdots \\ \vdots & \ddots & \\ \frac{\partial \mathcal{L}_n}{\partial W_{1M}} & & \frac{\partial \mathcal{L}_n}{\partial W_{DM}} \end{bmatrix}$$

$$\Rightarrow \quad \frac{\partial \mathcal{L}_n}{\partial W} = W^\top \left( \frac{\partial \mathcal{L}_n}{\partial C} + \frac{\partial \mathcal{L}_n}{\partial C}^\top \right)$$

## Probabilistic PCA

$\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}$ is symmetric, the matrix form of the derivatives are:

$$\frac{\partial}{\partial \mathbf{C}}(x_n - \mu)^\top \mathbf{C}^{-1}(x_n - \mu) = -\mathbf{C}^{-1}(x_n - \mu)(x_n - \mu)^\top \mathbf{C}^{-1}$$

$$\frac{\partial}{\partial \mathbf{C}} \log |\mathbf{C}| = \mathbf{C}^{-1}$$

Notice that both derivatives are symmetric matrices:

$$\frac{\partial \mathcal{L}_n}{\partial \mathbf{W}} = \mathbf{W}^\top \left( \frac{\partial \mathcal{L}_n}{\partial \mathbf{C}} + \frac{\partial \mathcal{L}_n}{\partial \mathbf{C}}^\top \right) = 2\mathbf{W}^\top \frac{\partial \mathcal{L}_n}{\partial \mathbf{C}}$$

Derivative of $\mathcal{L}$ w.r.t. $\mathbf{W}$:

$$\frac{\partial}{\partial \mathbf{W}} \log \mathcal{N}(x_n; \mu, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})$$

$$= 2\mathbf{W}^\top \frac{\partial}{\partial \mathbf{C}} \left( -\frac{1}{2}(x_n - \mu)^\top \mathbf{C}^{-1}(x_n - \mu) - \frac{1}{2} \log |\mathbf{C}| \right)$$

$$= \mathbf{W}^\top \left( \mathbf{C}^{-1}(x_n - \mu)(x_n - \mu)^\top \mathbf{C}^{-1} - \mathbf{C}^{-1} \right).$$

## Probabilistic PCA

Derivative of $\mathcal{L}$ w.r.t. $\mathbf{W}$ with $\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}$:

$$\frac{\partial}{\partial \mathbf{W}} \log \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I})$$

$$= 2\mathbf{W}^\top \frac{\partial}{\partial \mathbf{C}} \left( -\frac{1}{2}(\boldsymbol{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1}(\boldsymbol{x}_n - \boldsymbol{\mu}) - \frac{1}{2}\log|\mathbf{C}| \right)$$

$$= \mathbf{W}^\top \left( \mathbf{C}^{-1}(\boldsymbol{x}_n - \boldsymbol{\mu})(\boldsymbol{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1} - \mathbf{C}^{-1} \right).$$

$$\Rightarrow \quad \left( \frac{\partial \mathcal{L}}{\partial \mathbf{W}} \right)^\top = \mathbf{C}^{-1}(\underbrace{\frac{1}{N}\sum_{n=1}^{N}(\boldsymbol{x}_n - \boldsymbol{\mu})(\boldsymbol{x}_n - \boldsymbol{\mu})^\top}_{:=\mathbf{S}\text{ , covariance when }\boldsymbol{\mu}=\boldsymbol{\mu}^*}\mathbf{C}^{-1} - \mathbf{I})\mathbf{W}$$

Setting $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{W}^*$ satisfies $\mathbf{S}(\mathbf{W}^*(\mathbf{W}^*)^\top + \sigma^2\mathbf{I})^{-1}\mathbf{W}^* = \mathbf{W}^*$

# Probabilistic PCA

Setting $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{W}^*$ satisfies $\mathbf{S}(\mathbf{W}^*(\mathbf{W}^*)^\top + \sigma^2 \mathbf{I})^{-1}\mathbf{W}^* = \mathbf{W}^*$

Possible solutions for the fixed points:

1. $\mathbf{W}^* = \mathbf{0}$ (then $p_{\theta*}(x|z) = \mathcal{N}(x; \mu^*, \sigma^2 \mathbf{I})$, not interesting)

2. Lets write down the SVD of $\mathbf{W}^*$ and assume $\mathbf{W}^* = \mathbf{U}\Sigma\mathbf{V}^\top$,
   $\mathbf{U} \in \mathbb{R}^{D \times D}, \Sigma \in \mathbb{R}^{D \times M}, \mathbf{V} \in \mathbb{R}^{M \times M}$

$$
\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots \\ 0 & \ddots & \\ \vdots & & \sigma_M \\ & & 0 \\ & & \vdots \\ & & 0 \end{bmatrix} \Rightarrow \Sigma\Sigma^\top + \sigma^2 \mathbf{I} = \begin{bmatrix} \sigma_1^2 + \sigma^2 & 0 & \cdots & & & \\ 0 & \ddots & & & & \\ \vdots & & \sigma_M^2 + \sigma^2 & & & \\ & & & \sigma^2 & & \\ & & & & \ddots & \\ & & & & & \sigma^2 \end{bmatrix}
$$

## Probabilistic PCA

Setting $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{W}^*$ satisfies $\mathbf{S}(\mathbf{W}^*(\mathbf{W}^*)^\top + \sigma^2\mathbf{I})^{-1}\mathbf{W}^* = \mathbf{W}^*$

Possible solutions for the fixed points:

1. $\mathbf{W}^* = \mathbf{0}$ (then $p_{\theta*}(x|z) = \mathcal{N}(x; \mu^*, \sigma^2\mathbf{I})$, not interesting)

2. Lets write down the SVD of $\mathbf{W}^*$ and assume $\mathbf{W}^* = \mathbf{U}\Sigma\mathbf{V}^\top$,
   $\mathbf{U} \in \mathbb{R}^{D \times D}, \Sigma \in \mathbb{R}^{D \times M}, \mathbf{V} \in \mathbb{R}^{M \times M}$

$$\mathbf{S}(\mathbf{U}\Sigma\Sigma^\top\mathbf{U}^\top + \sigma^2\mathbf{I})^{-1}\mathbf{U}\Sigma\mathbf{V}^\top = \mathbf{U}\Sigma\mathbf{V}^\top$$
$$\Rightarrow \quad \mathbf{S}(\mathbf{U}(\Sigma\Sigma^\top + \sigma^2\mathbf{I})\mathbf{U}^\top)^{-1}\mathbf{U}\Sigma\mathbf{V}^\top = \mathbf{U}\Sigma\mathbf{V}^\top$$
$$\Rightarrow \quad \mathbf{S}\mathbf{U}(\Sigma\Sigma^\top + \sigma^2\mathbf{I})^{-1}\mathbf{U}^\top\mathbf{U}\Sigma\mathbf{V}^\top = \mathbf{U}\Sigma\mathbf{V}^\top$$
$$\Rightarrow \quad \mathbf{S}\mathbf{U}(\Sigma\Sigma^\top + \sigma^2\mathbf{I})^{-1}\Sigma\mathbf{V}^\top = \mathbf{U}\Sigma\mathbf{V}^\top$$
$$\Rightarrow \quad \mathbf{S}\mathbf{U}(\Sigma\Sigma^\top + \sigma^2\mathbf{I})^{-1}\Sigma = \mathbf{U}\Sigma$$

# Probabilistic PCA

Setting $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{W}^*$ satisfies $\mathbf{S}(\mathbf{W}^*(\mathbf{W}^*)^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^* = \mathbf{W}^*$

Possible solutions for the fixed points:

1. $\mathbf{W}^* = \mathbf{0}$ (then $p_{\boldsymbol{\theta}*}(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}^*, \sigma^2 \mathbf{I})$, not interesting)

2. Lets write down the SVD of $\mathbf{W}^*$ and assume $\mathbf{W}^* = \mathbf{U}\Sigma\mathbf{V}^\top$,
   $\mathbf{U} \in \mathbb{R}^{D \times D}, \Sigma \in \mathbb{R}^{D \times M}, \mathbf{V} \in \mathbb{R}^{M \times M}$

$$\mathbf{S}\mathbf{U} \begin{bmatrix} (\sigma_1^2 + \sigma^2)^{-1}\sigma_1 & 0 & \cdots & \\ 0 & \ddots & & \\ \vdots & & (\sigma_M^2 + \sigma^2)^{-1}\sigma_M & \\ & & 0 & \\ & & \vdots & \\ & & 0 & \end{bmatrix} = \mathbf{U} \begin{bmatrix} \sigma_1 & 0 & \cdots \\ 0 & \ddots & \\ \vdots & & \sigma_M \\ & & 0 \\ & & \vdots \\ & & 0 \end{bmatrix}$$

# Probabilistic PCA

Setting $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{W}^*$ satisfies $\mathbf{S}(\mathbf{W}^*(\mathbf{W}^*)^\top + \sigma^2 \mathbf{I})^{-1}\mathbf{W}^* = \mathbf{W}^*$

Possible solutions for the fixed points:

1. $\mathbf{W}^* = \mathbf{0}$ (then $p_{\boldsymbol{\theta}*}(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}^*, \sigma^2 \mathbf{I})$, not interesting)

2. Lets write down the SVD of $\mathbf{W}^*$ and assume $\mathbf{W}^* = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$,
   $\mathbf{U} \in \mathbb{R}^{D \times D}, \boldsymbol{\Sigma} \in \mathbb{R}^{D \times M}, \mathbf{V} \in \mathbb{R}^{M \times M}$

$$\mathbf{S}\mathbf{U}\begin{bmatrix} 1 & 0 & \cdots \\ 0 & \ddots & \\ \vdots & & 1 \\ & & 0 \\ & & \vdots \\ & & 0 \end{bmatrix} = \mathbf{U}\begin{bmatrix} \sigma_1^2 + \sigma^2 & 0 & \cdots \\ & 0 & \ddots \\ \vdots & & \sigma_M^2 + \sigma^2 \\ & & 0 \\ & & \vdots \\ & & 0 \end{bmatrix}$$

Write $\mathbf{U} := (\boldsymbol{u}_1, ..., \boldsymbol{u}_D)$:

$$(\mathbf{S}\boldsymbol{u}_1, ..., \mathbf{S}\boldsymbol{u}_M) = ((\sigma_1^2 + \sigma^2)\boldsymbol{u}_1, ..., (\sigma_M^2 + \sigma^2)\boldsymbol{u}_M)$$

$\Rightarrow$ the first $M$ columns of $\mathbf{U}$ contain eigenvectors of $\mathbf{S}$!

## Probabilistic PCA

Setting $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{W}^*$ satisfies $\mathbf{S}(\mathbf{W}^*(\mathbf{W}^*)^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^* = \mathbf{W}^*$

Possible solutions for the fixed points:

1. $\mathbf{W}^* = \mathbf{0}$ (then $p_{\boldsymbol{\theta}*}(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}^*, \sigma^2 \mathbf{I})$, not interesting)

2. Lets write down the SVD of $\mathbf{W}^*$ and assume $\mathbf{W}^* = \mathbf{U}\Sigma\mathbf{V}^\top$,
   $\mathbf{U} \in \mathbb{R}^{D \times D}, \Sigma \in \mathbb{R}^{D \times M}, \mathbf{V} \in \mathbb{R}^{M \times M}$
   $\Rightarrow \quad \mathbf{S}\mathbf{U}(\Sigma\Sigma^\top + \sigma^2 \mathbf{I})^{-1}\Sigma = \mathbf{U}\Sigma$
   Then given $\mathbf{S} = \mathbf{Q}\Lambda\mathbf{Q}^\top, \mathbf{Q} = (\boldsymbol{q}_1, ..., \boldsymbol{q}_D), \lambda_1 \geqslant ... \geqslant \lambda_D \geqslant 0$,

   $$\mathbf{U} := (\boldsymbol{u}_1, ..., \boldsymbol{u}_D), \boldsymbol{u}_m = \boldsymbol{q}_{i_m}, 1 \leqslant i_m \leqslant D, m = 1, ..., M$$

   ($\mathbf{U}$ can contain any other columns for $\boldsymbol{u}_{M+1}$ to $\boldsymbol{u}_D$)

# Probabilistic PCA

Setting $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{W}^*$ satisfies $\mathbf{S}(\mathbf{W}^*(\mathbf{W}^*)^\top + \sigma^2 \mathbf{I})^{-1}\mathbf{W}^* = \mathbf{W}^*$
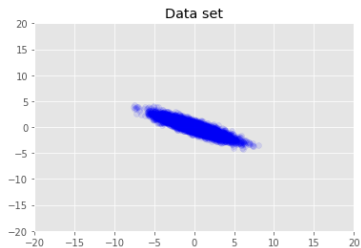
Possible solutions for the fixed points:

1. $\mathbf{W}^* = \mathbf{0}$ (then $p_{\boldsymbol{\theta}*}(x|z) = \mathcal{N}(x; \boldsymbol{\mu}^*, \sigma^2 \mathbf{I})$, not interesting)

2. Lets write down the SVD of $\mathbf{W}^*$ and assume $\mathbf{W}^* = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$,
   $\mathbf{U} \in \mathbb{R}^{D \times D}, \boldsymbol{\Sigma} \in \mathbb{R}^{D \times M}, \mathbf{V} \in \mathbb{R}^{M \times M}$
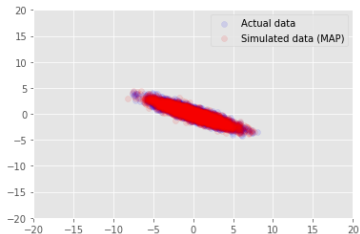   Then given $\mathbf{S} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top, \lambda_1 \geqslant ... \geqslant \lambda_D \geqslant 0$

   - **Exercise:** For $m = 1, ..., M$, $\Sigma_{mm} = \sqrt{\lambda_{i_m} - \sigma^2}$ if $u_m = q_{i_m}$
   - **Exercise:** Global maximum: $u_m = q_m$ for $mi = 1, ..., M$
     $\Rightarrow$ picking the $M$ principal components (like PCA)

# Probabilistic PCA



Dataset



Generate data with Prob. PCA

https://www.tensorflow.org/probability/examples/Probabilistic_PCA

# Extensions of Probabilistic PCA

Probabilistic PCA: make a latent variable model as follows:

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

$$p_{\theta}(x|z) = \mathcal{N}(x; \mathbf{W}z + \mu, \sigma^2 \mathbf{I})$$

From Probabilistic PCA to other interesting generative models:

- **Factor analysis**: change conditional output covariance from $\sigma^2 \mathbf{I}$ to $\Psi$ (a learnable diagonal matrix)

- **Generator for a VAE**: change conditional output mean from $\mathbf{W}z + \mu$ to $\mu_{\theta}(z)$ (See Deep Learning course next term)

- **Training**: (variational) expectation maximisation (See Probabilistic Inference course next term)

# Summary

Probabilistic PCA

- One of the simplest generative model (linear generator)

- Optimal solution closely related to PCA

One more exercise for you if you have time:
Derive the posterior $p_{\theta^*}(z|x)$ using the optimal $\theta^* = \{\mu^*, \mathbf{W}^*\}$