

# Principal Component Analysis

CoE197M/EE298M (Foundations of Machine Learning)

Rowel Atienza, Ph.D.

[rowel@eee.upd.edu.ph](mailto:rowel@eee.upd.edu.ph)

*Reference:* "Mathematics for Machine Learning". Copyright 2020 by Marc Peter Deisenroth, A. Aldo Faisal, and Cheng Soon Ong. Published by Cambridge University Press.

# Motivations

Raw data representations are over-complete

Dimensionality reduction reduces the footprint of data without losing useful important information

# Problem Statement

Consider a dataset  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ,  $\mathbf{x}_n \in \mathbb{R}^D$  with mean  $\mathbf{0}$  and data covariance matrix:

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

There exists a low-dimensional compression representation (code) of  $\mathbf{x}_n$ :

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M$$

The projection matrix:

$$\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$$

With orthonormal basis  $\mathbf{b}_i^T \mathbf{b}_j = 0$  with  $i \neq j$  and  $\mathbf{b}_i^T \mathbf{b}_i = 1$

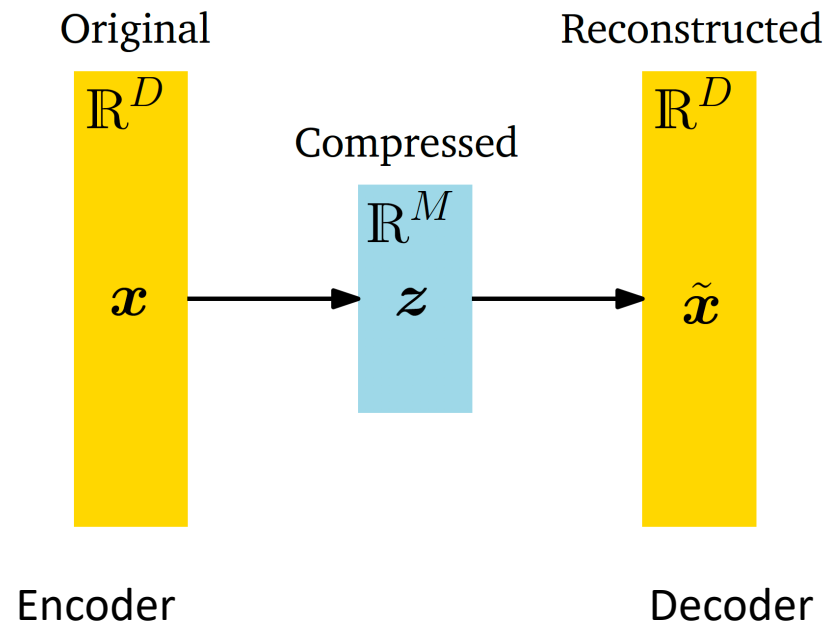
We project  $\mathbf{x}_n$  into a low-dimensional subspace  $U \subseteq \mathbb{R}^D$  with  $\dim(U) = M < D$

The projected data is  $\tilde{\mathbf{x}}_n$  with  $\mathbf{z}_n$  as the coordinates on basis  $\mathbf{B}$

# Dimensionality Reduction

The objective is to find  $\tilde{\mathbf{x}}_n \in \mathbb{R}^D$  or  $\mathbf{z}_n = [z_{1n}, \dots, z_{Mn}]^T \in \mathbb{R}^{M \times 1}$  and basis  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$  that minimizes the loss due to compression

Example loss: squared reconstruction loss  $\|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$



# Finding Projective Coordinates

# Projective Perspective

Assume ONB  $B = (\mathbf{b}_1, \dots, \mathbf{b}_D) \in \mathbb{R}^D$

$$\mathbf{x} = \sum_{i=1}^M k_i \mathbf{b}_i + \sum_{j=M+1}^D k_j \mathbf{b}_j$$



$$\tilde{\mathbf{x}}_n \in U \subseteq \mathbb{R}^D$$

where  $k_i \in \mathbb{R}$

# Objective Function

Minimize the average Euclidean reconstruction error:

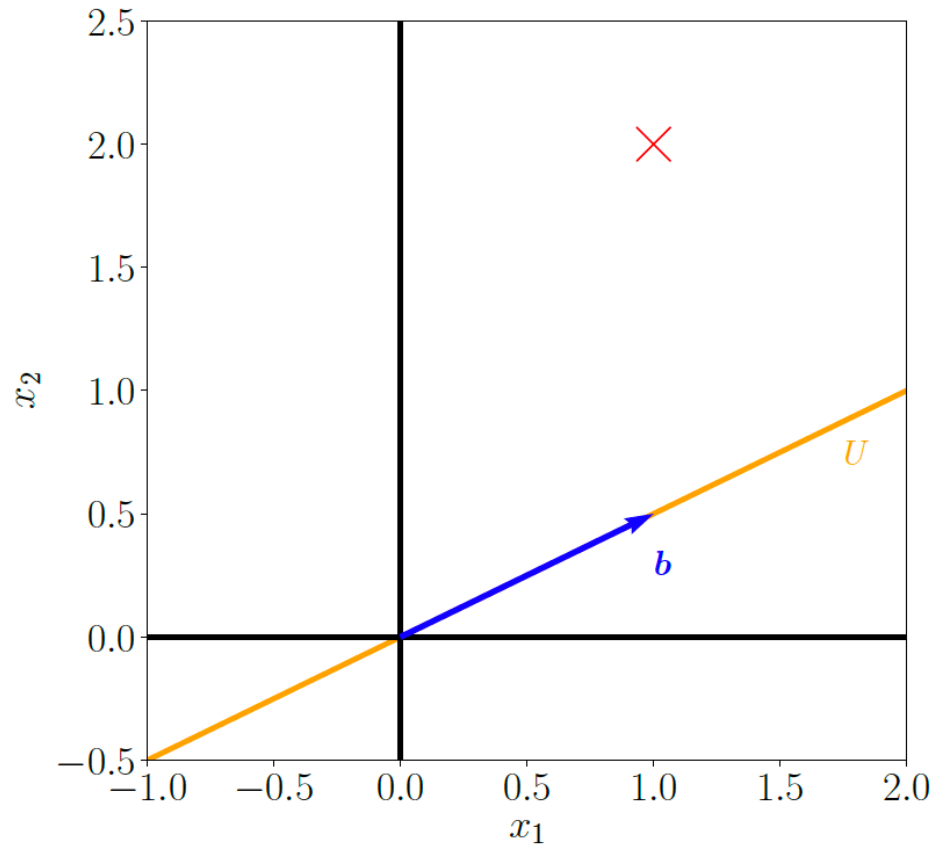
$$J_M := \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

Where

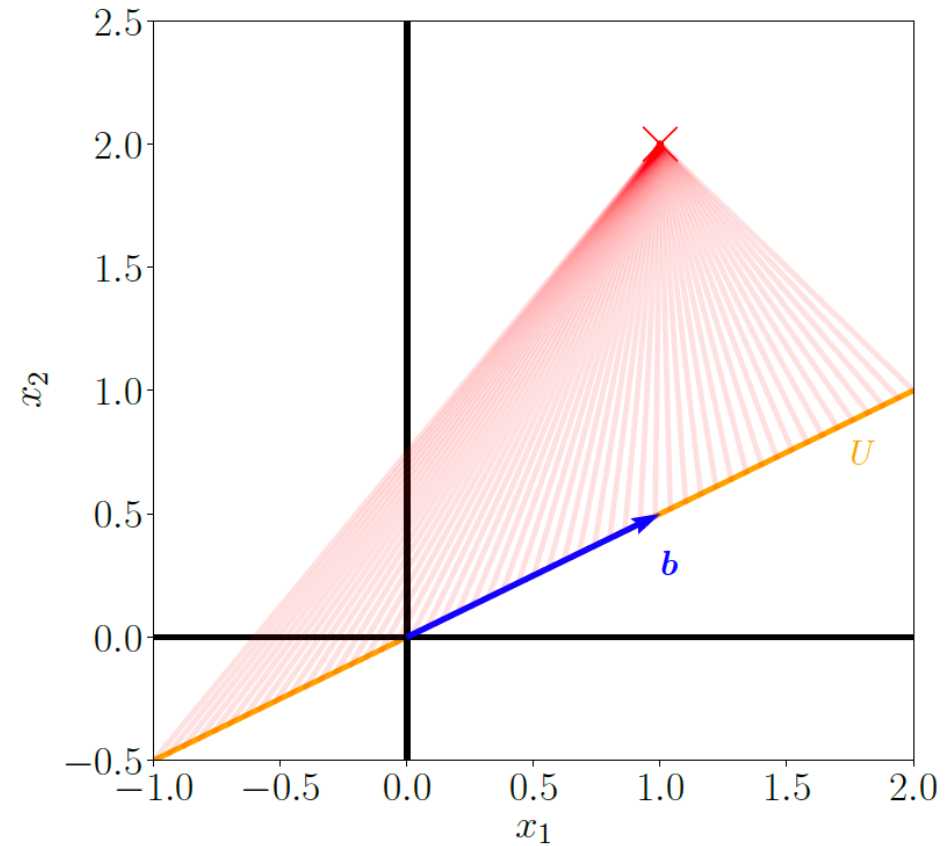
$$\tilde{\mathbf{x}}_n = \sum_{m=1}^M z_{mn} \mathbf{b}_m = \mathbf{B} \mathbf{z}_n \in U \subseteq \mathbb{R}^D$$



# Visualizing Distance Minimization

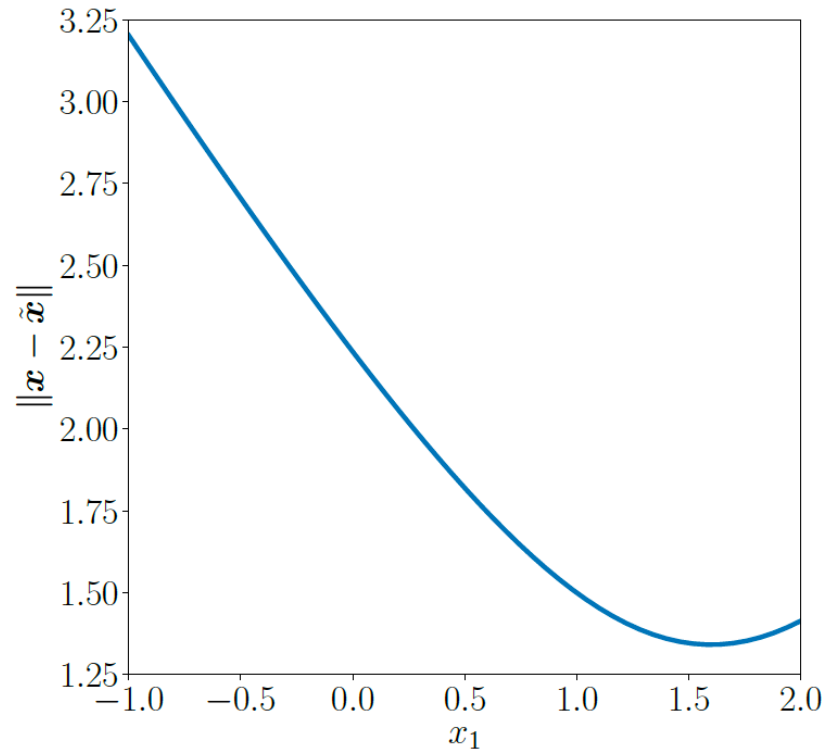


(a) Setting.

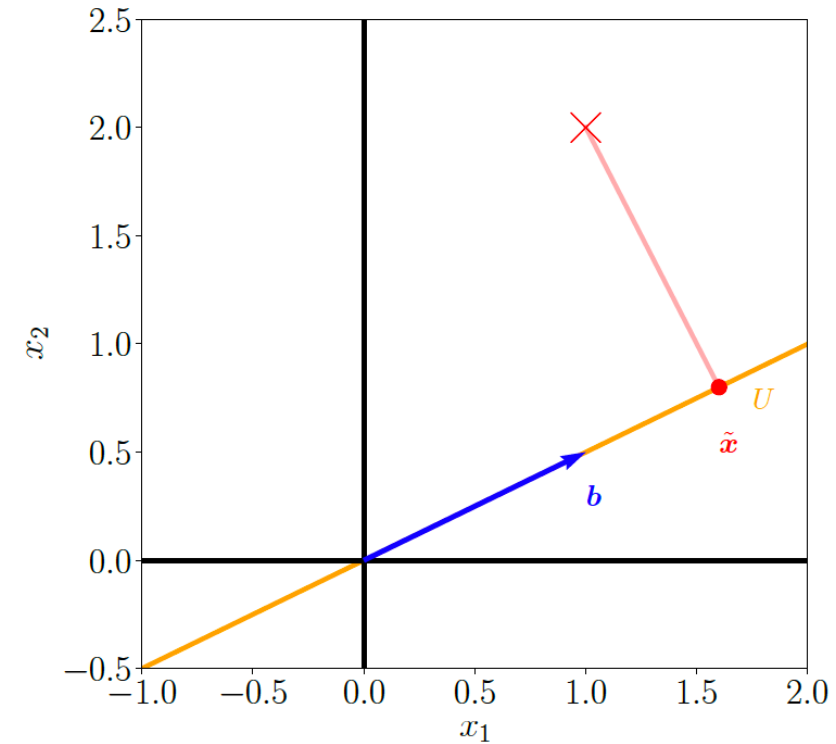


(b) Differences  $x - \tilde{x}_i$  for 50 different  $\tilde{x}_i$  are shown by the red lines.

# Visualizing Distance Minimization



(a) Distances  $\|x - \tilde{x}\|$  for some  $\tilde{x} = z_1 \mathbf{b} \in U = \text{span}[\mathbf{b}]$ ; see panel (b) for the setting.



(b) The vector  $\tilde{x}$  that minimizes the distance in panel (a) is its orthogonal projection onto  $U$ . The coordinate of the projection  $\tilde{x}$  with respect to the basis vector  $\mathbf{b}$  that spans  $U$  is the factor we need to scale  $\mathbf{b}$  in order to “reach”  $\tilde{x}$ .

# Optimal Coordinates

Find coordinates of  $\mathbf{z}$  of  $\tilde{\mathbf{x}}_n$  for  $n = 1, \dots, N$

Assume, ONB  $[\mathbf{b}_1, \dots, \mathbf{b}_M]$  of  $U \subseteq \mathbb{R}^D$

Such that  $\|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$  is minimized

Assume we are given dataset:  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  where  $\mathbf{x}_n \in \mathbb{R}^D$  and  $\mathbb{E}[\mathcal{X}] = 0$  or all data are zero centered

# Zero Centering

Assuming the mean  $\mathbb{E}[\mathcal{X}] = \boldsymbol{\mu}$

$$\boldsymbol{x}_n = \boldsymbol{x}_n - \boldsymbol{\mu}$$

# Optimal Coordinates

$$\frac{\partial J_M}{\partial z_{in}} = \frac{\partial J_M}{\partial \tilde{\mathbf{x}}_n} \frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}}$$

$$\frac{\partial J_M}{\partial \tilde{\mathbf{x}}_n} = -\frac{2}{N} (\mathbf{x}_n - \tilde{\mathbf{x}}_n)^T \in \mathbb{R}^{1 \times D}$$

$$\frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}} = \frac{\partial}{\partial z_{in}} \left( \sum_{m=1}^M z_{mn} \mathbf{b}_m \right) = \mathbf{b}_i \in \mathbb{R}^{D \times 1}$$

# Optimal Projection $z_{in}$

$$\frac{\partial J_M}{\partial z_{in}} = -\frac{2}{N} (\mathbf{x}_n - \tilde{\mathbf{x}}_n)^T \mathbf{b}_i = -\frac{2}{N} \left( \mathbf{x}_n - \sum_{m=1}^M z_{mn} \mathbf{b}_m \right)^T \mathbf{b}_i$$

$$\frac{\partial J_M}{\partial z_{in}} = -\frac{2}{N} (\mathbf{x}_n^T \mathbf{b}_i - z_{in} \mathbf{b}_i^T \mathbf{b}_i)$$

$$\frac{\partial J_M}{\partial z_{in}} = -\frac{2}{N} (\mathbf{x}_n^T \mathbf{b}_i - z_{in}) = 0$$

$$z_{in} = \mathbf{x}_n^T \mathbf{b}_i = \mathbf{b}_i^T \mathbf{x}_n$$

# Observations

The optimal projection  $\tilde{\mathbf{x}}_n$  of  $\mathbf{x}_n$  is an orthogonal projection

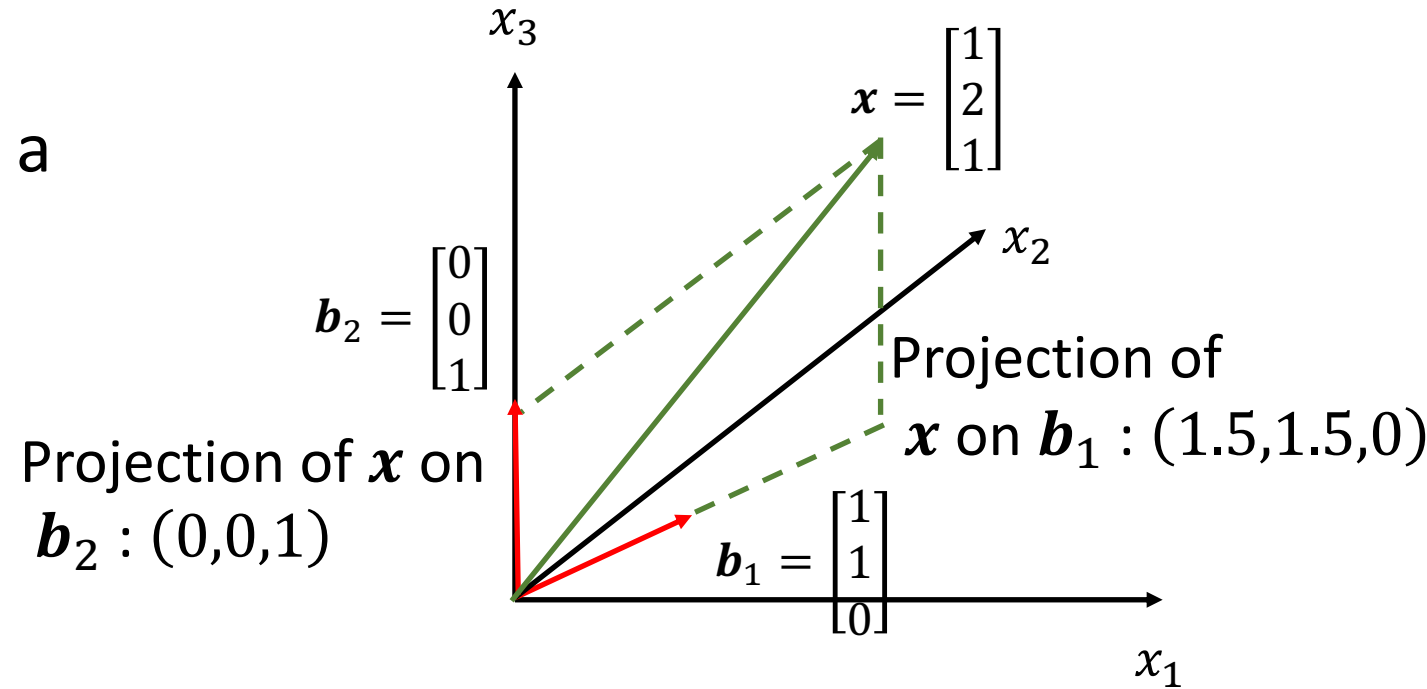
The coordinates of  $\tilde{\mathbf{x}}_n$  with respect to  $[\mathbf{b}_1, \dots, \mathbf{b}_M]$  are the coordinates of the orthogonal projection of  $\tilde{\mathbf{x}}_n$  on the principal subspace

An orthogonal projection is the best linear mapping given the objective

The coordinates  $z_{in}$  for  $i = 1, \dots, m$  must be the same as  $k_{in}$  for  $i = 1, \dots, m$

# PCA on 2-dim

The coordinates of  $\mathbf{x}$  on  $\mathbf{b}_1$  has a length  $\mathbf{z}_1 = \sqrt{5}$ .



Projection of  $\mathbf{x}$  on  $\mathbf{b}_1$  has length  $\sqrt{1 + 2^2} = \sqrt{5}$ .

This is split to:

$$x_1 = \sqrt{5} \cos \frac{\pi}{4} = 1.5 \text{ and } x_2 = \sqrt{5} \sin \frac{\pi}{4} = 1.5$$



# Orthogonal Projection with ONB

Recall projection of a vector  $\mathbf{x} \in \mathbb{R}^n$  onto  $U$  that is closest to  $\mathbf{x}$  is  $\pi_U(\mathbf{x})$  with a basis vector  $\mathbf{b} \in \mathbb{R}^n$

$$\pi_U(\mathbf{x}) = \lambda \mathbf{b} = \mathbf{b} \lambda = \mathbf{b} \frac{\mathbf{b}^T \mathbf{x}}{\|\mathbf{b}\|^2} = \frac{\mathbf{b} \mathbf{b}^T}{\|\mathbf{b}\|^2} \mathbf{x} = \mathbf{P}_\pi \mathbf{x}$$

Then:

$$\tilde{\mathbf{x}}_n = \mathbf{b}_i \underbrace{\mathbf{b}_i^T \mathbf{x}_n}_{z_{in}}$$

$$z_{in} = \mathbf{b}_i^T \mathbf{x}_n$$

# General ONB

Assume, ONB  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M]$  of  $U \subseteq \mathbb{R}^D$

$$\tilde{\mathbf{x}}_n = \mathbf{B}\mathbf{B}^T \mathbf{x}_n$$

$\mathbf{B}^T \mathbf{x}_n$  is the projection of  $\mathbf{x}_n$  on ONB

Note that  $\tilde{\mathbf{x}}_n \in \mathbb{R}^D$  but our coordinates  $[z_1, \dots, z_M]$  with respect to basis vectors  $[\mathbf{b}_1, \dots, \mathbf{b}_M]$  is of dimensions  $M < D$

The other coordinates  $[z_{M+1}, \dots, z_D]$  with respect to basis vectors  $[\mathbf{b}_{M+1}, \dots, \mathbf{b}_D]$  have zero values

Finding ONB

Find the Basis  $\mathbf{b}_1, \dots, \mathbf{b}_M$

$$\tilde{\mathbf{x}}_n = \sum_{m=1}^M z_{mn} \mathbf{b}_m$$

$$\tilde{\mathbf{x}}_n = \sum_{m=1}^M (\mathbf{x}_n^T \mathbf{b}_m) \mathbf{b}_m$$

$$\tilde{\mathbf{x}}_n = \left( \sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^T \right) \mathbf{x}_n$$

$$\mathbf{x}_n = \sum_{m=1}^M z_m \mathbf{b}_m + \sum_{j=M+1}^D z_j \mathbf{b}_j$$

$$\mathbf{x}_n = \left( \sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^T \right) \mathbf{x}_n + \left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^T \right) \mathbf{x}_n$$

Therefore,

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^T \right) \mathbf{x}_n = \sum_{j=M+1}^D (\mathbf{x}_n^T \mathbf{b}_j) \mathbf{b}_j$$

Observation on  $\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{j=M+1}^D (\mathbf{x}_n^T \mathbf{b}_j) \mathbf{b}_j$

The difference is exactly the projection of the data point on the orthogonal complement of the principal subspace

# Maximum Variance

Project to low-dimensional subspace while maximizing variance to retain as much information as possible

# Finding the 1<sup>st</sup> Basis Vector $\mathbf{b}_1 \in \mathbb{R}^D$

Assuming i.i.d., maximize the variance of the first coordinate  $z_1$  of  $\mathbf{z} \in \mathbb{R}^M$ :

$$V_1 = \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2$$

Where

$$z_{1n} = \mathbf{b}_1^T \mathbf{x}_n$$



$$V_1 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^T \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \mathbf{b}_1^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{b}_1$$

$$V_1 = \mathbf{b}_1^T \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_1 = \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1$$

Where  $\mathbf{S}$  is the data covariance matrix defined earlier.

We restrict  $\|\mathbf{b}_1\|^2 = 1$  so that the variance comes from  $\mathbf{S}$  only

# Direction of Maximum Variance

$$\max_{\mathbf{b}_1} \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1$$

Subject to:  $\|\mathbf{b}_1\|^2 = 1$

The Lagrange:

$$\mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 + \lambda(1 - \mathbf{b}_1^T \mathbf{b}_1)$$

The partial derivative:

$$\frac{d\mathcal{L}}{d\mathbf{b}_1} = 2\mathbf{b}_1^T \mathbf{S} - 2\lambda \mathbf{b}_1^T$$

$$\frac{d\mathcal{L}}{d\lambda} = (1 - \mathbf{b}_1^T \mathbf{b}_1)$$

$$2\mathbf{b}_1^T \mathbf{S} - 2\lambda \mathbf{b}_1^T = 0 \text{ and } 1 - \mathbf{b}_1^T \mathbf{b}_1 = 0:$$

$$\mathbf{b}_1^T \mathbf{S} = \lambda \mathbf{b}_1^T \text{ or } \mathbf{S} \mathbf{b}_1 = \lambda \mathbf{b}_1$$
$$\mathbf{b}_1^T \mathbf{b}_1 = 1$$

Rewriting:

$$V_1 = \mathbf{b}_1^T \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_1 = \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 = \lambda \mathbf{b}_1^T \mathbf{b}_1 = \lambda$$

The variance is equal to the eigenvalue.

To maximize the variance, we choose the eigenvector as the basis vector with the maximum eigenvalue.

This eigenvector  $\mathbf{b}_1$  is called first principal component

# What we have so far...

We have one basis vector  $\mathbf{b}_1$  which is the eigenvector corresponding to the largest eigenvalue of  $\mathbf{S}$

Problem: We need  $m - 1$  more basis vectors  $\mathbf{b}_2, \dots, \mathbf{b}_m$

## Finding the 2<sup>nd</sup> Basis Vector $\mathbf{b}_2 \in \mathbb{R}^D$

Subtract the effect of the first principal component  $\mathbf{b}_1$  from the data:

$$\hat{\mathbf{x}}_n = \mathbf{x}_n - \tilde{\mathbf{x}}_1 = \mathbf{x}_n - \mathbf{b}_1 \mathbf{b}_1^T \mathbf{x}_n$$

Then we can use the same argument:

$$V_2 = \mathbf{b}_2^T \left( \frac{1}{N} \sum_{n=1}^N \hat{\mathbf{x}}_n \hat{\mathbf{x}}_n^T \right) \mathbf{b}_2 = \mathbf{b}_2^T \hat{\mathbf{S}} \mathbf{b}_2$$

The maximum variance is achieved at

$$\hat{\mathbf{S}} \mathbf{b}_2 = \lambda \mathbf{b}_2$$

Finding the  $M^{th}$  Basis Vector  $\mathbf{b}_M \in \mathbb{R}^D$

Subtract the effect of the first  $M - 1$  principal components  $\mathbf{b}_1, \dots, \mathbf{b}_{M-1}$  from the data:

$$\hat{\mathbf{x}}_n = \mathbf{x}_n - \left( \sum_{m=1}^{M-1} \mathbf{b}_m \mathbf{b}_m^T \right) \mathbf{x}_n$$

Then we can use the same argument:

$$V_m = \mathbf{b}_m^T \left( \frac{1}{N} \sum_{n=1}^N \hat{\mathbf{x}}_n \hat{\mathbf{x}}_n^T \right) \mathbf{b}_m = \mathbf{b}_m^T \hat{\mathbf{S}} \mathbf{b}_m$$

The maximum variance is achieved at

$$\hat{\mathbf{S}} \mathbf{b}_m = \lambda \mathbf{b}_m$$

# Recall: Properties of Symmetric Matrix

*Theorem (Spectral Theorem):* If  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric, there exists an **orthonormal basis** of vector space  $V$  from the eigenvectors of  $\mathbf{A}$  and each eigenvalue is real.

*Theorem:* The eigenvectors  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  of matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with distinct eigenvalues  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  are linearly independent

*Theorem (SPSD):* For a given matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , we can always obtain a symmetric positive semi-definite matrix  $\mathbf{S} \in \mathbb{R}^{n \times n}$ :  $\mathbf{S} = \mathbf{A}^T \mathbf{A}$

If  $\text{rank}(\mathbf{A}) = n$ , then  $\mathbf{S}$  is a symmetric positive definite (SPD) matrix

# Eigenvalues/Eigenvectors

$$\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, \mathbf{x}_n \in \mathbb{R}^D :$$

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T = \frac{1}{N} \mathbf{X} \mathbf{X}^T$$

$$\mathbf{X} = [\mathbf{x}_1 \quad \cdots \quad \mathbf{x}_N] \in \mathbb{R}^{D \times N}$$



# Typical Procedure to Obtain Eigenvalues/Eigenvectors

Perform Eigendecomposition on  $\mathbf{S} = \mathbf{B}\mathbf{D}\mathbf{B}^{-1}$

$\mathbf{D}$  are eigenvalues

$\mathbf{B}$  are eigenvectors

Perform SVD on  $\mathbf{X}$

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$
$$\mathbf{S} = \frac{1}{N}\mathbf{X}\mathbf{X}^T = \frac{1}{N}\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{V}\mathbf{\Sigma}^T\mathbf{U}^T = \frac{1}{N}\mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^T\mathbf{U}^T$$

$\mathbf{\Sigma}$  are eigenvalues

$\mathbf{U}$  are eigenvectors

# Low-Rank Approximations of $\mathbf{X}$

Consider the SVD of  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$

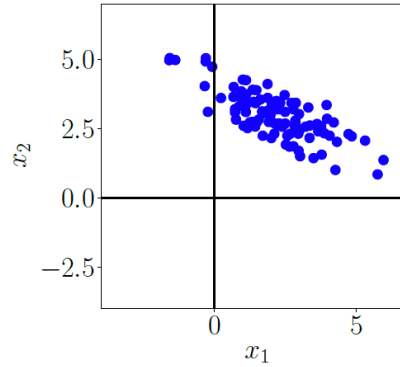
A low-rank approximation of  $\mathbf{X}$  using the  $M$  largest eigenvalues:

$$\mathbf{X} = \mathbf{U}_M \mathbf{\Sigma}_M \mathbf{V}_M^T \in \mathbb{R}^{D \times N}$$

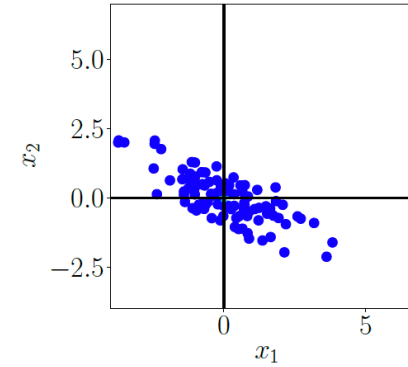
# PCA Algorithm

1. Mean Subtraction :  $\mathbf{x}_n = \mathbf{x}_n - \boldsymbol{\mu}$  where  $\boldsymbol{\mu} = [u_1 \quad \cdots \quad u_d]^T$ ,  $d$  is the data dimension (eg  $d = 3$  for RGB image)
2. Standardization by dividing the data by standard deviation:  $\mathbf{x}_n = \frac{\mathbf{x}_n - \boldsymbol{\mu}}{\boldsymbol{\sigma}}$ , where  $\boldsymbol{\sigma} = [\sigma_1 \quad \cdots \quad \sigma_d]^T$
3. Eigendecomposition of covariance matrix  $\mathbf{S} = \mathbf{B}\mathbf{D}\mathbf{B}^{-1}$
4. Projection:  $\tilde{\mathbf{x}}_n = \mathbf{B}_M \mathbf{B}_M^T \mathbf{x}_n$  where the coordinates with respect to the  $M$  principal basis vectors subspace:  $\tilde{\mathbf{z}}_n = \mathbf{B}_M^T \mathbf{x}_n$
5. Backprojection:  $\tilde{\mathbf{x}}_n = \tilde{\mathbf{x}}_n \boldsymbol{\sigma} + \boldsymbol{\mu}$

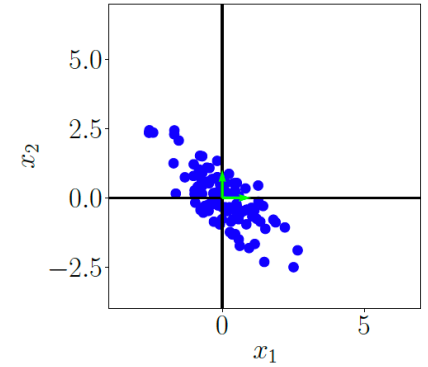
# PCA Algorithm



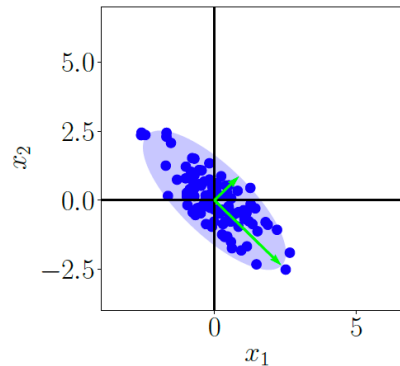
(a) Original dataset.



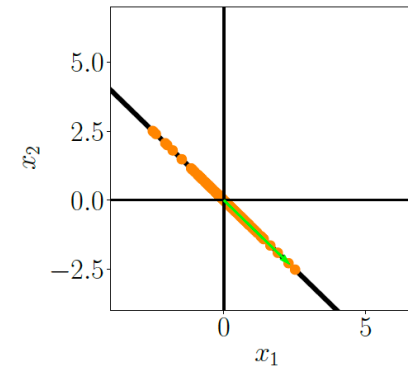
(b) Step 1: Centering by subtracting the mean from each data point.



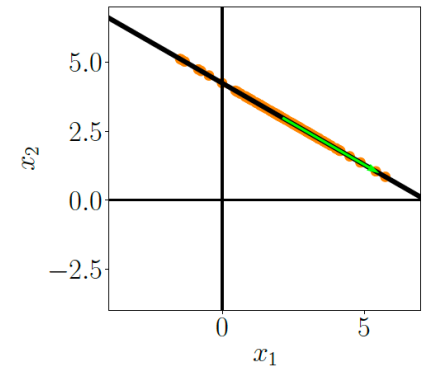
(c) Step 2: Dividing by the standard deviation to make the data unit free. Data has variance 1 along each axis.



(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).

# Probabilistic Modelling

# Probabilistic Model

Typical probabilistic model:

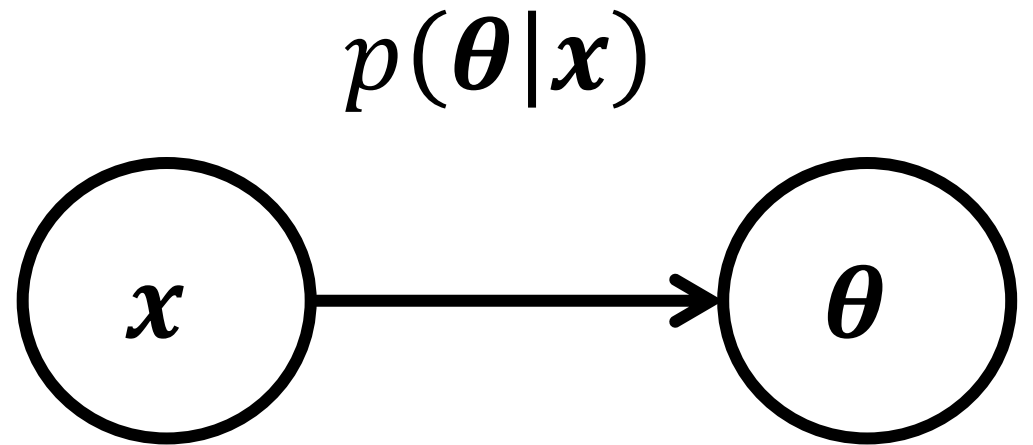
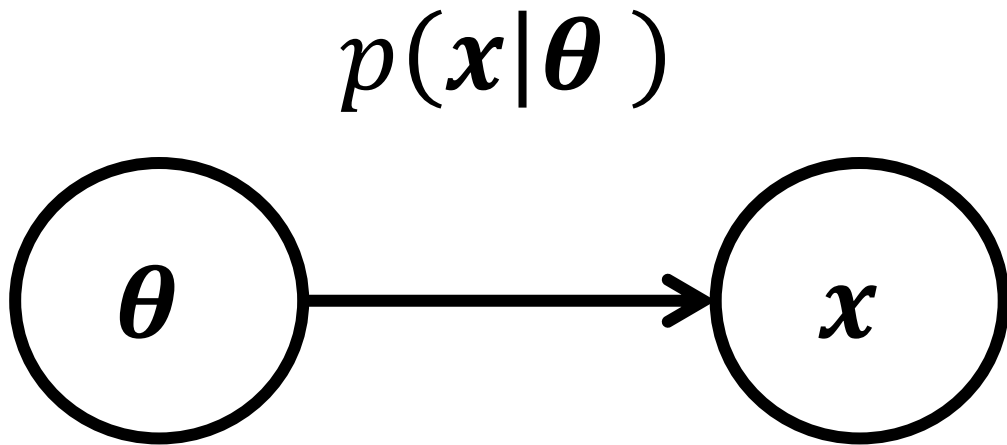
$$\begin{array}{c} \text{Joint probability} \qquad \text{likelihood} \quad \text{prior} \\ \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \\ p(\boldsymbol{x}, \boldsymbol{\theta}) = p(\boldsymbol{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \end{array}$$

Marginal:  $p(\boldsymbol{x}) = \int p(\boldsymbol{x}, \boldsymbol{\theta}) d\boldsymbol{\theta}$

Posterior:  $p(\boldsymbol{\theta}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}, \boldsymbol{\theta})}{p(\boldsymbol{x})}$

# Probabilistic Model DGM

$$p(\boldsymbol{x}, \boldsymbol{\theta}) = p(\boldsymbol{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\boldsymbol{x})p(\boldsymbol{x})$$



# Latent Variable Model

An intermediate latent variable  $\mathbf{z}$  is introduced as part of the model

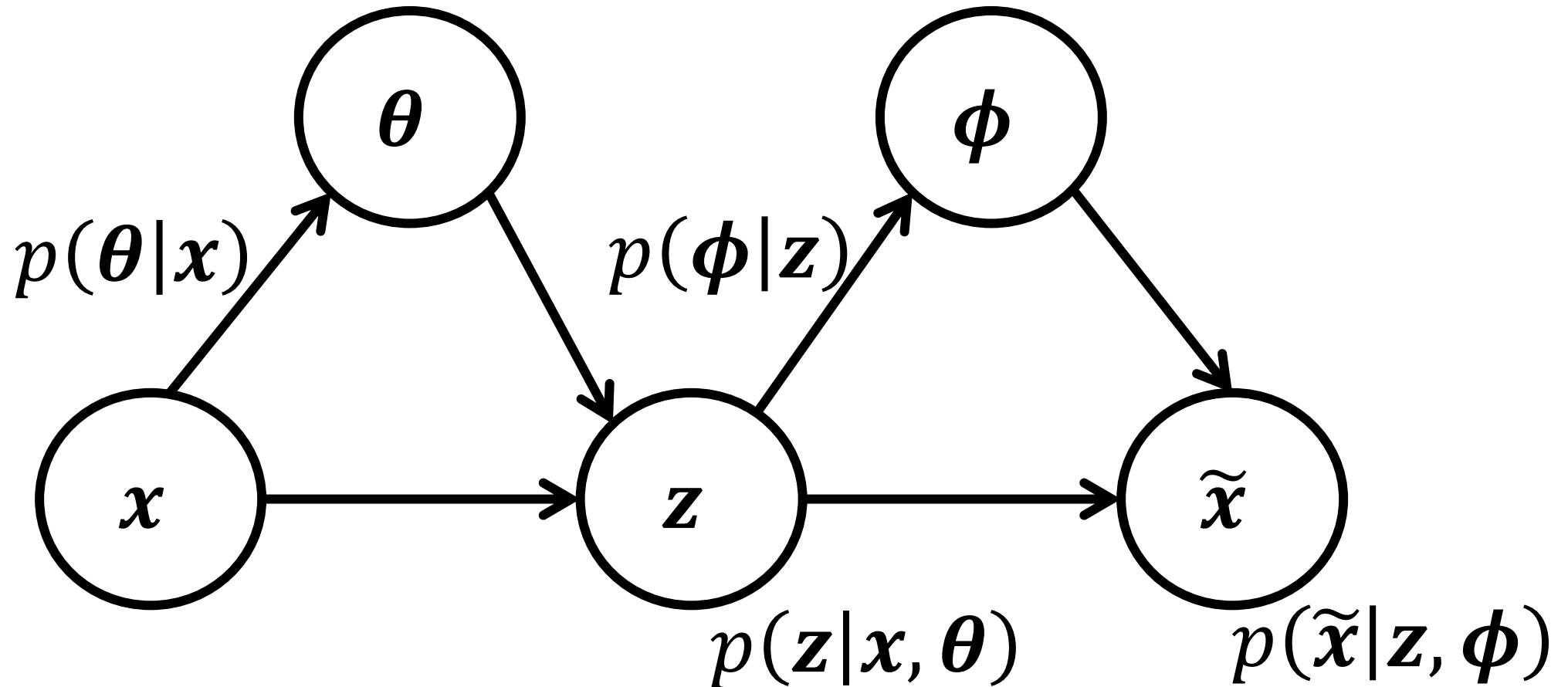
The latent variable  $\mathbf{z}$  is not a model parameter

The latent variable  $\mathbf{z}$  describes both the data distribution  $p(\mathbf{x})$ , thus the data generating process  $p(\tilde{\mathbf{x}}|\mathbf{z}, \boldsymbol{\phi})$  where  $\boldsymbol{\phi}$  represents the model parameters



# Latent Variable Model DGM

$$p(\mathbf{x}, \tilde{\mathbf{x}}, \boldsymbol{\theta}, \boldsymbol{\phi}) = p(\mathbf{x})p(\boldsymbol{\theta}|\mathbf{x})p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\phi}|\mathbf{z})p(\tilde{\mathbf{x}}|\mathbf{z}, \boldsymbol{\phi})$$



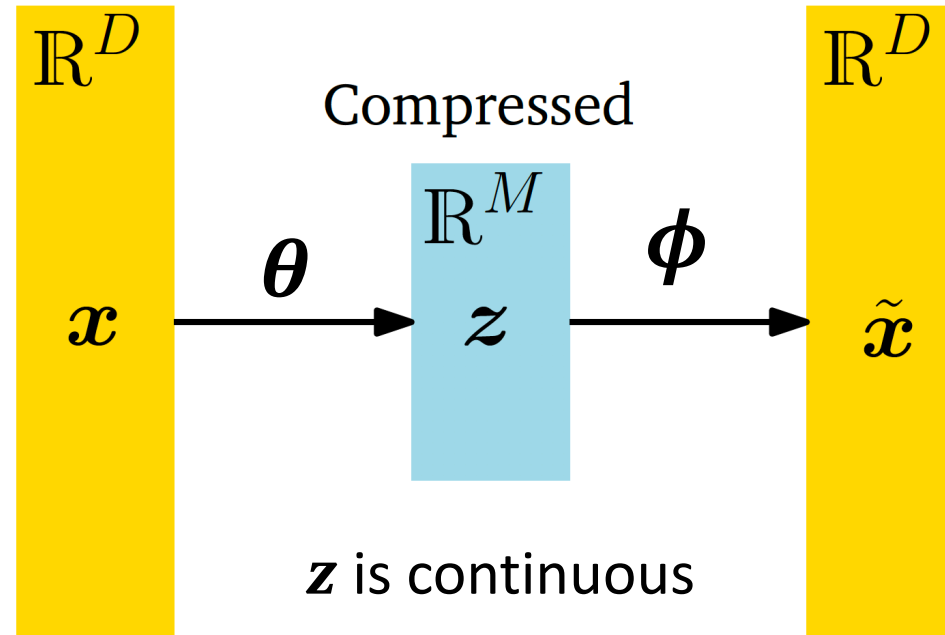
# Latent Variable Model of PCA

Encoder:  $p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$

Decoder:  $p(\tilde{\mathbf{x}}|\mathbf{z}, \boldsymbol{\phi})$

Original

Reconstructed



# Probabilistic PCA (PPCA)

Can deal with noise

Can use Bayesian interpretation

Can use PCA decoder as a generator

Can generate new data points from the generator

Can extend to mixture of PCA

Can treat PCA as a special case

etc

# PPCA

If  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$  and linear relationship between latent variable and observed data  $\mathbf{x}$ ,

$$\mathbf{x} = \mathbf{B}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} \in \mathbb{R}^D$$

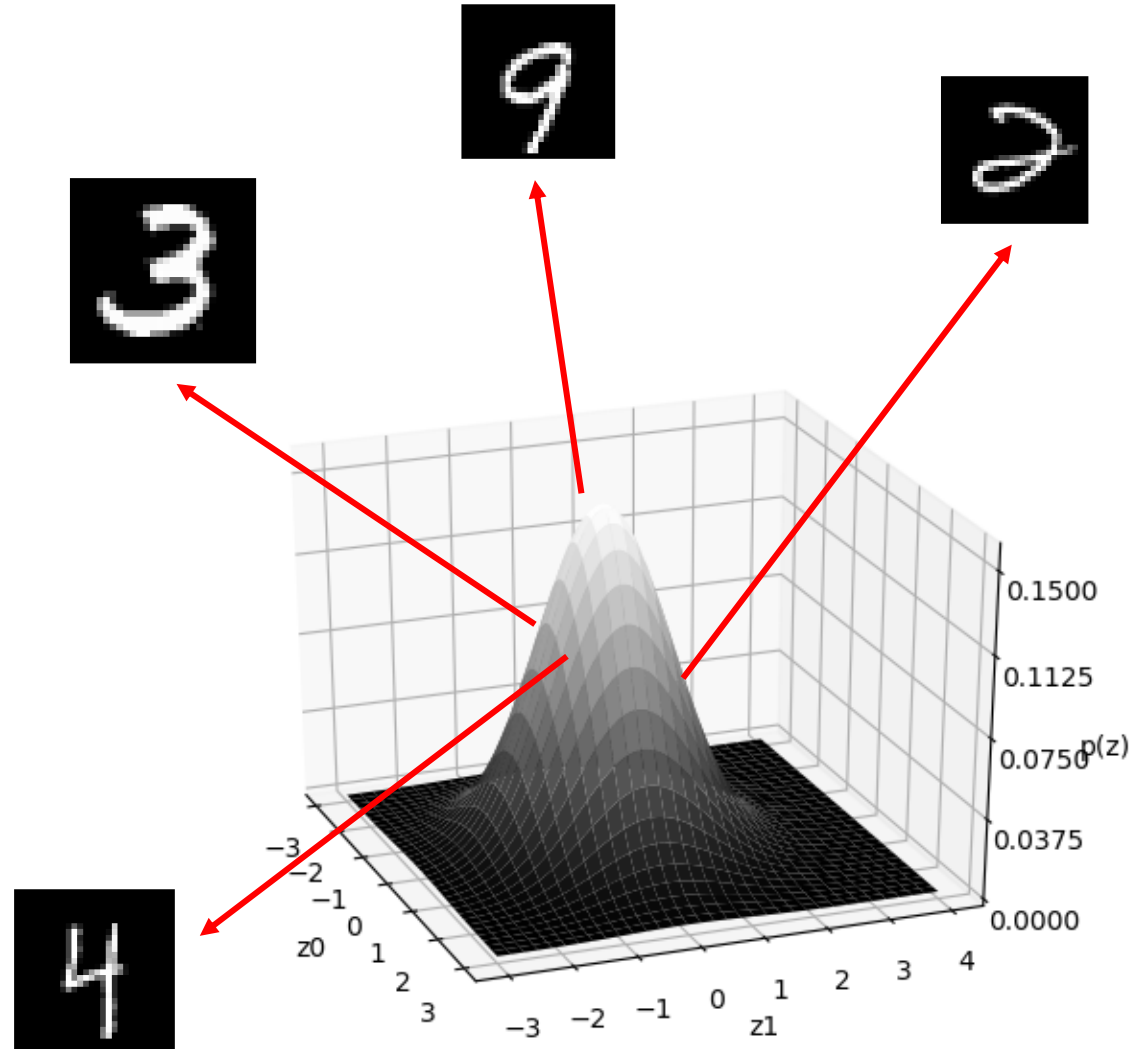
$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ ,  $\mathbf{B} \in \mathbb{R}^{D \times M}$ ,  $\boldsymbol{\mu} \in \mathbb{R}^D$ :

$$p(\mathbf{x} | \mathbf{B}, \mathbf{z}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\mathbf{x} | \mathbf{B}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

# Generative Model

$$\mathbf{z}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{x} | \mathbf{z}_n \sim \mathcal{N}(\mathbf{x} | \mathbf{B}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$



# Generative Model

$$\begin{aligned} p(\boldsymbol{x}|\boldsymbol{B}, \boldsymbol{\mu}, \sigma^2) &= \int p(\boldsymbol{x}|\boldsymbol{B}, \boldsymbol{\mu}, \boldsymbol{z}, \sigma^2) p(\boldsymbol{z}) d\boldsymbol{z} \\ &= \int \mathcal{N}(\boldsymbol{x}|\boldsymbol{B}\boldsymbol{z} + \boldsymbol{\mu}, \sigma^2 \boldsymbol{I}) \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}) d\boldsymbol{z} \\ &= \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{B}\boldsymbol{B}^T + \sigma^2 \boldsymbol{I}) \end{aligned}$$