AIMS | African Institute for Mathematical Sciences | RWANDA

Imperial College London

# Principal Component Analysis

**Marc Deisenroth**

Quantum Leap Africa
African Institute for Mathematical
Sciences, Rwanda

Department of Computing
Imperial College London

@mpd37
mdeisenroth@aimsammi.org

October 4, 2018

# References

- Bishop: Pattern Recognition and Machine Learning, Chapter 12
- Deisenroth et al.: Mathematics for Machine Learning, Chapter 10
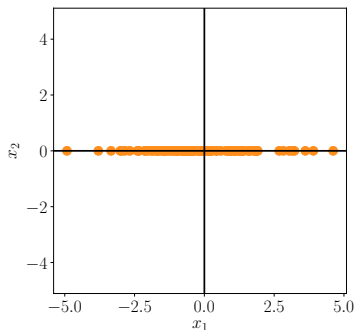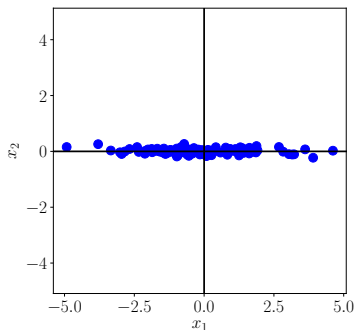  (https://mml-book.com)

# Overview

## Introduction

# High-Dimensional Data



- Real-world data is often high dimensional
- **Challenges:**
  - Difficult to analyze
  - Difficult to visualize
  - Difficult to interpret

# Properties of High-dimensional Data



- ‣ Many dimensions are unnecessary
- ‣ Data often lives on a low-dimensional manifold

⮞ Dimensionality reduction finds the relevant dimensions.

# Background: Coordinate Representations

Consider $\mathbb{R}^2$ with the canonical basis $\boldsymbol{e}_1 = [1,0]^\top$, $\boldsymbol{e}_2 = [0,1]^\top$.

$$\boldsymbol{x} = \begin{bmatrix} 5 \\ 3 \end{bmatrix} = 5\boldsymbol{e}_1 + 3\boldsymbol{e}_2 \qquad \text{Linear combination of basis vectors}$$

▸ **Coordinates** of $\boldsymbol{x}$ w.r.t. $(\boldsymbol{e}_1, \boldsymbol{e}_2)$: $[5, 3]$

# Background: Coordinate Representations

Consider $\mathbb{R}^2$ with the canonical basis $e_1 = [1,0]^\top$, $e_2 = [0,1]^\top$.

$$x = \begin{bmatrix} 5 \\ 3 \end{bmatrix} = 5e_1 + 3e_2 \qquad \text{Linear combination of basis vectors}$$

- **Coordinates** of $x$ w.r.t. $(e_1, e_2)$: $[5, 3]$

Consider the vectors of the form

$$\tilde{x} = \begin{bmatrix} 0 \\ z \end{bmatrix} \in \mathbb{R}^2, \quad z \in \mathbb{R}$$

▶▶ Write them as $0e_1 + ze_2$.

- Only remember/store the **coordinate/code** $z$ of the $e_2$ vector
  ▶▶ **Compression**
- Set of $\tilde{x}$ vectors forms a vector subspace $U \subseteq \mathbb{R}^2$ with $\dim(U) = 1$
  because $U = \text{span}[e_2]$.

# Overview

# PCA Setting



- Dataset $\mathcal{X} := \{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\}$, $\boldsymbol{x}_n \in \mathbb{R}^D$
- Data matrix $\boldsymbol{X} := [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N] \in \mathbb{R}^{D \times N}$ ▶▶ Often $N \times D$ matrix

# PCA Setting



original $\mathbb{R}^D$  —  code $\mathbb{R}^M$  —  compressed $\mathbb{R}^D$

$x$ → $z$ → $\tilde{x}$

Encoder  Decoder

- Dataset $\mathcal{X} := \{x_1, \ldots, x_N\}$, $x_n \in \mathbb{R}^D$
- Data matrix $X := [x_1, \ldots, x_N] \in \mathbb{R}^{D \times N}$ ▶▶ Often $N \times D$ matrix
- Without loss of generality: $\mathbb{E}[\mathcal{X}] = 0$ ▶▶ Centered data
  ▶▶ Data covariance matrix
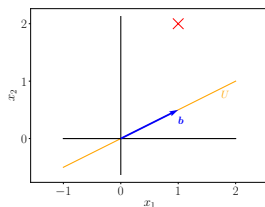
# PCA Setting



▸ Dataset $\mathcal{X} := \{x_1, \ldots, x_N\}$, $x_n \in \mathbb{R}^D$

▸ Data matrix $X := [x_1, \ldots, x_N] \in \mathbb{R}^{D \times N}$ ▶▶ Often $N \times D$ matrix

▸ Without loss of generality: $\mathbb{E}[\mathcal{X}] = \mathbf{0}$ ▶▶ Centered data
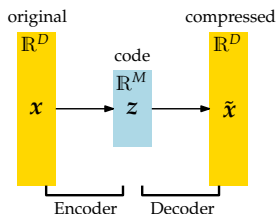  ▶▶ Data covariance matrix $S = \frac{1}{N} X X^\top \in \mathbb{R}^{D \times D}$

▸ Linear relationships between latent code $z$ and data $x$:

$$z = B^\top x, \quad \tilde{x} = Bz$$

▸ $B = [b_1, \ldots, b_M] \in \mathbb{R}^{D \times M}$ is an orthogonal matrix

# Low-Dimensional Embedding



- Find an $M$-dimensional subspace $U \subset \mathbb{R}^D$ onto which we project the data
- $\tilde{x} = \pi_U(x)$ is the projection of $x$ onto $U$
- Find projections $\tilde{x}$ that are as similar to $x$ as possible
  - ▶▶ Find basis vectors $b_1, \ldots, b_M$
- Compression loss incurs if $M \ll D$

# Overview

# PCA Idea: Maximum Variance



- Project $D$-dimensional data $x$ onto an $M$-dimensional subspace
  that retains as much information as possible
  ▶▶ Data compression

# PCA Idea: Maximum Variance



- ▸ Project $D$-dimensional data $x$ onto an $M$-dimensional subspace that retains as much information as possible
  - ▸▸ Data compression
- ▸ Informally: information = diversity = variance
  - ▸▸ **Maximize variance in projected space** (Hotelling 1933)

# PCA Objective: Maximum Variance

‣ Linear relationships:

$$z = B^\top x, \quad \tilde{x} = Bz$$

‣ $B = [b_1, \ldots, b_M] \in \mathbb{R}^{D \times M}$ is an orthogonal matrix
‣ Columns of $B$ are an ONB of an $M$-dimensional subspace of $\mathbb{R}^D$

# PCA Objective: Maximum Variance

- Linear relationships:

$$z = B^\top x, \quad \tilde{x} = Bz$$

- $B = [b_1, \ldots, b_M] \in \mathbb{R}^{D \times M}$ is an orthogonal matrix
- Columns of $B$ are an ONB of an $M$-dimensional subspace of $\mathbb{R}^D$
- Find $B = [b_1, \ldots, b_M]$ so that the variance in the projected space is maximized

$$\max_{b_1, \ldots, b_M} \mathbb{V}[z] = \max_{b_1, \ldots, b_M} \mathbb{V}[B^\top x]$$

$$\text{s.t. } \|b_1\| = 1 = \ldots = \|b_M\|$$

▶▶ Constrained optimization problem

# Direction with Maximal Variance (1)

‣ Maximize variance of first coordinate of $z \in \mathbb{R}^M$:

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^{N} z_{n1}^2$$

▶▶ Empirical variance of the training dataset

# Direction with Maximal Variance (1)

‣ Maximize variance of first coordinate of $z \in \mathbb{R}^M$:

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^{N} z_{n1}^2$$

▶ Empirical variance of the training dataset

‣ First coordinate of $z_n$ is

$$z_{n1} = b_1^\top x_n$$

▶ Coordinate of orthogonal projection of $x_n$ onto span$[b_1]$
(1-dimensional subspace spanned by $b_1$)

# Direction with Maximal Variance (1)

▸ Maximize variance of first coordinate of $z \in \mathbb{R}^M$:

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^{N} z_{n1}^2$$

▶ Empirical variance of the training dataset

▸ First coordinate of $z_n$ is

$$z_{n1} = b_1^\top x_n$$

▶ Coordinate of orthogonal projection of $x_n$ onto span$[b_1]$
(1-dimensional subspace spanned by $b_1$)

$$\mathbb{V}[z_1] =$$

# Direction with Maximal Variance (1)

- Maximize variance of first coordinate of $z \in \mathbb{R}^M$:

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^{N} z_{n1}^2$$

  ▶▶ Empirical variance of the training dataset

- First coordinate of $z_n$ is

$$z_{n1} = b_1^\top x_n$$

  ▶▶ Coordinate of orthogonal projection of $x_n$ onto span$[b_1]$
  (1-dimensional subspace spanned by $b_1$)

$$\mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^{N} (b_1^\top x_n)^2 = \frac{1}{N} \sum_{n=1}^{N} b_1^\top x_n x_n^\top b_1$$
$$= b_1^\top \left( \frac{1}{N} \sum_{n=1}^{N} x_n x_n^\top \right) b_1 = b_1^\top S b_1$$

# Direction with Maximal Variance (2)

- Maximize variance

$$\max_{\boldsymbol{b}_1, \|\boldsymbol{b}_1\|^2 = 1} \mathbb{V}[z_1] = \max_{\boldsymbol{b}_1, \|\boldsymbol{b}_1\|^2 = 1} \boldsymbol{b}_1^\top \boldsymbol{S} \boldsymbol{b}_1$$

# Direction with Maximal Variance (2)

▸ Maximize variance

$$\max_{\boldsymbol{b}_1, \|\boldsymbol{b}_1\|^2=1} \mathbb{V}[z_1] = \max_{\boldsymbol{b}_1, \|\boldsymbol{b}_1\|^2=1} \boldsymbol{b}_1^\top \boldsymbol{S} \boldsymbol{b}_1$$

▸ Lagrangian:

$$L(\boldsymbol{b}_1, \lambda) = \boldsymbol{b}_1^\top \boldsymbol{S} \boldsymbol{b}_1 + \lambda_1 (1 - \boldsymbol{b}_1^\top \boldsymbol{b}_1)$$

**Discuss with your neighbors and find $\lambda_1$ and $b_1$**

# Direction with Maximal Variance (2)

▸ Maximize variance

$$\max_{\boldsymbol{b}_1, \|\boldsymbol{b}_1\|^2=1} \mathbb{V}[z_1] = \max_{\boldsymbol{b}_1, \|\boldsymbol{b}_1\|^2=1} \boldsymbol{b}_1^\top \boldsymbol{S} \boldsymbol{b}_1$$

▸ Lagrangian:

$$L(\boldsymbol{b}_1, \lambda) = \boldsymbol{b}_1^\top \boldsymbol{S} \boldsymbol{b}_1 + \lambda_1 (1 - \boldsymbol{b}_1^\top \boldsymbol{b}_1)$$

**Discuss with your neighbors and find $\lambda_1$ and $b_1$**

▸ Setting the gradients w.r.t. $\boldsymbol{b}_1$ and $\lambda_1$ to **0** yields

$$\boldsymbol{S} \boldsymbol{b}_1 = \lambda_1 \boldsymbol{b}_1$$
$$\boldsymbol{b}_1^\top \boldsymbol{b}_1 = 1$$

▸ $\boldsymbol{b}_1$ is an eigenvector of the data covariance matrix $\boldsymbol{S}$
▸ $\lambda_1$ is the corresponding eigenvalue

# Direction with Maximal Variance (3)

- $Sb_1 = \lambda_1 b_1$

$$\mathbb{V}[z_1] = b_1^\top S b_1 = \lambda_1 b_1^\top b_1 = \lambda_1$$

▶▶ Variance retained by first coordinate corresponds to
eigenvalue $\lambda_1$

# Direction with Maximal Variance (3)

- $Sb_1 = \lambda_1 b_1$

$$\mathbb{V}[z_1] = b_1^\top S b_1 = \lambda_1 b_1^\top b_1 = \lambda_1$$

▶▶ Variance retained by first coordinate corresponds to eigenvalue $\lambda_1$

▶▶ Choose eigenvector $b_1$ associated with the largest eigenvalue

# Direction with Maximal Variance (3)

▸ $\boldsymbol{S}\boldsymbol{b}_1 = \lambda_1\boldsymbol{b}_1$

$$\mathbb{V}[z_1] = \boldsymbol{b}_1^\top\boldsymbol{S}\boldsymbol{b}_1 = \lambda_1\boldsymbol{b}_1^\top\boldsymbol{b}_1 = \lambda_1$$

▶▶ Variance retained by first coordinate corresponds to eigenvalue $\lambda_1$

▶▶ Choose eigenvector $\boldsymbol{b}_1$ associated with the largest eigenvalue

▸ Projection:

▸ Coordinate:

## Direction with Maximal Variance

Maximizing the variance means to choose the direction $\boldsymbol{b}_1$ as the eigenvector of the data covariance matrix $\boldsymbol{S}$ that is associated with the largest eigenvalue $\lambda_1$ of $\boldsymbol{S}$.

# Direction with Maximal Variance (3)

▸ $Sb_1 = \lambda_1 b_1$

$$\mathbb{V}[z_1] = b_1^\top S b_1 = \lambda_1 b_1^\top b_1 = \lambda_1$$

▸▸ Variance retained by first coordinate corresponds to eigenvalue $\lambda_1$

▸▸ Choose eigenvector $b_1$ associated with the largest eigenvalue

▸ Projection: $\tilde{x}_n = b_1 b_1^\top x_n$

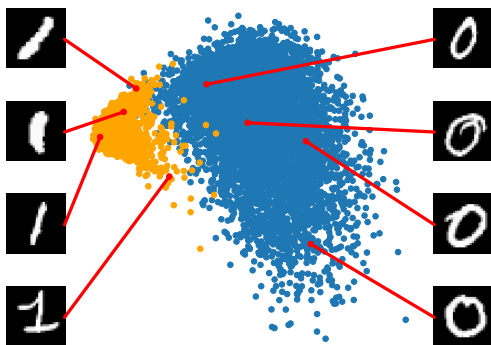▸ Coordinate: $z_{n1} = b_1^\top x_n$

---

### Direction with Maximal Variance

Maximizing the variance means to choose the direction $b_1$ as the eigenvector of the data covariance matrix $S$ that is associated with the largest eigenvalue $\lambda_1$ of $S$.

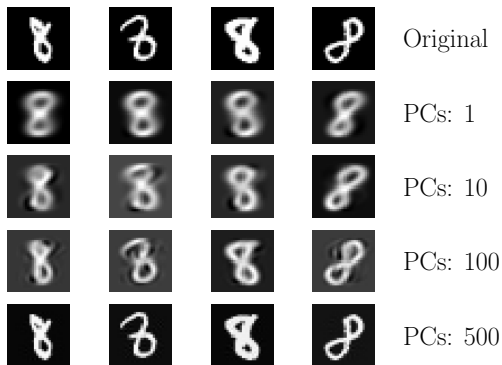# *M*-dimensional Subspace with Maximum Variance

### General Result

The *M*-dimensional subspace of $\mathbb{R}^D$ that retains the most variance is spanned by the *M* eigenvectors of the data covariance matrix $S$ that are associated with the *M* largest eigenvalues of $S$. (e.g., Bishop 2006)

# Example: MNIST Embedding (Training Set)



▸ Embedding of handwritten '0' and '1' digits (28 × 28 pixels) into
  a two-dimensional subspace, spanned by the first two principal
  components.

# Example: MNIST Reconstruction (Test Set)



Original

PCs: 1

PCs: 10

PCs: 100

PCs: 500

- ▸ Reconstructions of original digits as the number of principal components increases

# Overview

# Refresher: Orthogonal Projection onto Subspaces

- Basis $b_1, \ldots, b_M$ of a subspace $U \subset \mathbb{R}^D$
- Define $B = [b_1, ..., b_M] \in \mathbb{R}^{D \times M}$
- Project $x \in \mathbb{R}^D$ onto subspace $U$:

$$\pi_U(x) = \tilde{x} = B(B^\top B)^{-1} B^\top x$$

- If $b_1, \ldots, b_M$ form an orthonormal basis ($b_i^\top b_j = \delta_{ij}$), then the projection simplifies to

$$\tilde{x} = BB^\top x$$

# PCA Objective: Minimize Reconstruction Error



▸ Objective: Find orthogonal projection that minimizes the average squared projection/reconstruction error

$$J = \frac{1}{N} \sum_{n=1}^{N} \|\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n\|^2$$

where $\tilde{\boldsymbol{x}}_n = \pi_U(\boldsymbol{x}_n)$ is the projection of $\boldsymbol{x}_n$ onto $U$

# Derivation (1)

- Assume an orthonormal basis of $\mathbb{R}^D = \text{span}[\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D]$, such that $\boldsymbol{b}_i^\top \boldsymbol{b}_j = \delta_{ij}$

# Derivation (1)

- Assume an orthonormal basis of $\mathbb{R}^D = \text{span}[\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D]$, such that $\boldsymbol{b}_i^\top \boldsymbol{b}_j = \delta_{ij}$
- Every data point $\boldsymbol{x}$ can be written as a linear combination of the basis vectors:

$$\boldsymbol{x} = \sum_{d=1}^{D} \eta_d \boldsymbol{b}_d = \boldsymbol{B}\boldsymbol{\eta}, \quad \boldsymbol{B} = [\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D]$$

# Derivation (1)

- Assume an orthonormal basis of $\mathbb{R}^D = \text{span}[\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D]$, such that $\boldsymbol{b}_i^\top \boldsymbol{b}_j = \delta_{ij}$
- Every data point $\boldsymbol{x}$ can be written as a linear combination of the basis vectors:

$$\boldsymbol{x} = \sum_{d=1}^{D} \eta_d \boldsymbol{b}_d = \boldsymbol{B}\boldsymbol{\eta}, \quad \boldsymbol{B} = [\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D]$$

▶▶ Rotation of the standard coordinates to a new coordinate system defined by the basis $(\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D)$.

# Derivation (1)

- Assume an orthonormal basis of $\mathbb{R}^D = \mathrm{span}[\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D]$, such that $\boldsymbol{b}_i^\top \boldsymbol{b}_j = \delta_{ij}$
- Every data point $\boldsymbol{x}$ can be written as a linear combination of the basis vectors:

$$\boldsymbol{x} = \sum_{d=1}^{D} \eta_d \boldsymbol{b}_d = \boldsymbol{B}\boldsymbol{\eta}, \quad \boldsymbol{B} = [\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D]$$

▶▶ Rotation of the standard coordinates to a new coordinate system defined by the basis $(\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D)$.

▶▶ Original coordinates $x_d$ are replaced by $\eta_d, d = 1, \ldots, D$

# Derivation (1)

- Assume an orthonormal basis of $\mathbb{R}^D = \mathrm{span}[\boldsymbol{b}_1, \dots, \boldsymbol{b}_D]$, such that $\boldsymbol{b}_i^\top \boldsymbol{b}_j = \delta_{ij}$
- Every data point $\boldsymbol{x}$ can be written as a linear combination of the basis vectors:

$$\boldsymbol{x} = \sum_{d=1}^{D} \eta_d \boldsymbol{b}_d = \boldsymbol{B}\boldsymbol{\eta}, \quad \boldsymbol{B} = [\boldsymbol{b}_1, \dots, \boldsymbol{b}_D]$$

  ▶▶ Rotation of the standard coordinates to a new coordinate system defined by the basis $(\boldsymbol{b}_1, \dots, \boldsymbol{b}_D)$.

  ▶▶ Original coordinates $x_d$ are replaced by $\eta_d$, $d = 1, \dots, D$

- Obtain $\eta_d = \boldsymbol{x}^\top \boldsymbol{b}_d$, such that

$$\boldsymbol{x} = \sum_{d=1}^{D} (\boldsymbol{x}^\top \boldsymbol{b}_d) \boldsymbol{b}_d$$

# Derivation (2)

## Objective

Approximate

$$\boldsymbol{x} = \sum_{d=1}^{D} \eta_d \boldsymbol{b}_d \quad \text{with} \quad \tilde{\boldsymbol{x}} = \sum_{m=1}^{M} z_m \boldsymbol{b}_m$$

using $M \ll D$ many basis vectors

▸▸ **Projection** onto a lower-dimensional subspace

# Derivation (2)

## Objective

Approximate

$$\boldsymbol{x} = \sum_{d=1}^{D} \eta_d \boldsymbol{b}_d \quad \text{with} \quad \tilde{\boldsymbol{x}} = \sum_{m=1}^{M} z_m \boldsymbol{b}_m$$

using $M \ll D$ many basis vectors
▶▶ **Projection** onto a lower-dimensional subspace

# Derivation (3): Objective

$$\tilde{\boldsymbol{x}}_n = \underbrace{\sum_{m=1}^{M} z_{mn} \boldsymbol{b}_m}_{\text{lower-dim. subspace}}$$

# Derivation (3): Objective

$$\tilde{\boldsymbol{x}}_n = \underbrace{\sum_{m=1}^{M} z_{mn}\boldsymbol{b}_m}_{\text{lower-dim. subspace}}$$

▸ Choose coordinates $z_{mn}$ and basis vectors $\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D$ such that the average squared reconstruction error

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

is minimized

# Derivation (3): Objective

$$\tilde{\boldsymbol{x}}_n = \underbrace{\sum_{m=1}^{M} z_{mn} \boldsymbol{b}_m}_{\text{lower-dim. subspace}}$$

‣ Choose coordinates $z_{mn}$ and basis vectors $\boldsymbol{b}_1, \ldots, \boldsymbol{b}_D$ such that the average squared reconstruction error

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

is minimized

▶▶ Compute gradients of $J_M$ w.r.t. all variables, set to $\boldsymbol{0}$, solve

# Derivation (4): Optimal Coordinates

Necessary condition for optimum:

$$\frac{\partial J_M}{\partial z_{mn}} = 0 \quad \implies \quad z_{mn} = \boldsymbol{x}_n^\top \boldsymbol{b}_m, \qquad m = 1, \dots, M$$

- The optimal projection is the orthogonal projection
- The optimal coordinate $z_{mn}$ is the orthogonal projection of $\boldsymbol{x}_n$ onto the one-dimensional subspace spanned by $\boldsymbol{b}_m$
- $(\boldsymbol{b}_1, \dots, \boldsymbol{b}_D)$ is ONB ▶▶ span$[\boldsymbol{b}_{M+1}, \dots, \boldsymbol{b}_D]$ is orthogonal complement of principal subspace (span$[\boldsymbol{b}_1, \dots, \boldsymbol{b}_M]$)
- If

$$\boldsymbol{x}_n = \sum_{d=1}^{D} \eta_{dn} \boldsymbol{b}_d \quad \text{and} \quad \tilde{\boldsymbol{x}}_n = \sum_{m=1}^{M} z_{mn} \boldsymbol{b}_m$$

then $\eta_{mn} = z_{mn}$ for $m = 1, \dots, M$

▶▶ Minimum error is given by the orthogonal projection of $\boldsymbol{x}_n$ onto the principal subspace spanned by $\boldsymbol{b}_1, \dots, \boldsymbol{b}_M$
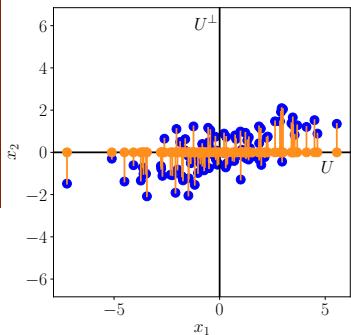
# Derivation (5): Displacement Vector



Approximation error only plays a role in dimensions $M + 1, \ldots, D$:

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

# Derivation (5): Displacement Vector



Approximation error only plays a role in dimensions $M + 1, \ldots, D$:

$$x_n - \tilde{x}_n = \sum_{j=M+1}^{D} (x_n^\top b_j) b_j$$

▶▶ Displacement vector $x_n - \tilde{x}_n$ lies in orthogonal complement $U^\perp$ of principal subspace $U$ (linear combination of the $b_j$ for $j = M + 1, \ldots, D$)

# Derivation (5)

From the previous slide:

$$\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n = \sum_{j=M+1}^{D} (\boldsymbol{x}_n^\top \boldsymbol{b}_j) \boldsymbol{b}_j$$

# Derivation (5)

From the previous slide:

$$\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n = \sum_{j=M+1}^{D} (\boldsymbol{x}_n^\top \boldsymbol{b}_j)\boldsymbol{b}_j$$

Let's compute our reconstruction error:

$$J_M = \frac{1}{N} \sum_{n=1}^{N} \|\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n\|^2 = \frac{1}{N} \sum_{n=1}^{N} (\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n)^\top (\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n)$$

# Derivation (5)

From the previous slide:

$$\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n = \sum_{j=M+1}^{D} (\boldsymbol{x}_n^\top \boldsymbol{b}_j) \boldsymbol{b}_j$$

Let's compute our reconstruction error:

$$J_M = \frac{1}{N} \sum_{n=1}^{N} \|\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n\|^2 = \frac{1}{N} \sum_{n=1}^{N} (\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n)^\top (\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \sum_{j=M+1}^{D} (\boldsymbol{x}_n^\top \boldsymbol{b}_j)^2$$

# Derivation (5)

From the previous slide:

$$x_n - \tilde{x}_n = \sum_{j=M+1}^{D} (x_n^\top b_j) b_j$$

Let's compute our reconstruction error:

$$J_M = \frac{1}{N} \sum_{n=1}^{N} \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^{N} (x_n - \tilde{x}_n)^\top (x_n - \tilde{x}_n)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \sum_{j=M+1}^{D} (x_n^\top b_j)^2$$

$$= \sum_{j=M+1}^{D} b_j^\top S b_j$$

where $S = \frac{1}{N} \sum_{n=1}^{N} x_n x_n^\top$ is the data covariance matrix

# Derivation (6)

- What remains: Minimize $J_M$ w.r.t. $\boldsymbol{b}_j$ under the constraint that the $\boldsymbol{b}_j$ form an orthonormal basis.

- Similar setting to maximum variance perspective: Instead of maximizing the variance in the principal subspace, we minimize the variance in the orthogonal complement of the principal subspace

- End up with **eigenvalue problem**:

$$\boldsymbol{S}\boldsymbol{b}_j = \lambda_j \boldsymbol{b}_j, \quad j = D+1, \dots, M$$

# Derivation (7)

‣ Find the eigenvectors $b_j$ of the data covariance matrix $S$

# Derivation (7)

- Find the eigenvectors $\boldsymbol{b}_j$ of the data covariance matrix $\boldsymbol{S}$
- Corresponding value of the squared reconstruction error:

$$J_M = \sum_{j=M+1}^{D} \lambda_j$$

i.e., the sum of the eigenvalues associated with eigenvectors not in the principle subspace

# Derivation (7)

- Find the eigenvectors $\boldsymbol{b}_j$ of the data covariance matrix $\boldsymbol{S}$
- Corresponding value of the squared reconstruction error:

$$J_M = \sum_{j=M+1}^{D} \lambda_j$$

  i.e., the sum of the eigenvalues associated with eigenvectors not in the principle subspace

- Minimizing $J_M$ requires us to choose the $M$ eigenvectors as the principle subspace that are associated with the $M$ largest eigenvalues.

# Geometric Interpretation



‣ Objective: Project $x$ onto an affine subspace $\mu + \mathrm{span}[b_1]$.

# Geometric Interpretation



- Shift scenario to the origin (affine subspace ⇝ vector subspace)

# Geometric Interpretation



- Shift $x$ as well (onto $x - \mu$).

# Geometric Interpretation



- Orthogonal projection of $x - \mu$ onto subspace spanned by $b_1$

# Geometric Interpretation



▸ Move projected point $\pi_{U_1}(\boldsymbol{x})$ back into original (affine) setting.

# Overview

# Key Steps of PCA

1. Compute the empirical mean $\mu$ of the data
2. Mean subtraction: Replace all data points $x_i$ with $\bar{x}_i = x_i - \mu$.
3. Standardization: Divide the data by its standard deviation in each dimension: $\hat{X}^{(d)} = \bar{X}/\sigma(X^{(d)})$ for $d = 1, \ldots, D$.
4. Eigendecomposition of the data covariance matrix: Compute the eigenvectors (orthonormal) and eigenvalues of the data covariance matrix $S$
5. Orthogonal projection: Choose the eigenvectors associated with the $M$ largest eigenvalues to be the basis of the principal subspace. Obtain $\tilde{X}$
6. Moving back to original data space: $\tilde{X}^{(d)} = \tilde{X}^{(d)}\sigma(X^{(d)}) + \mu_d$
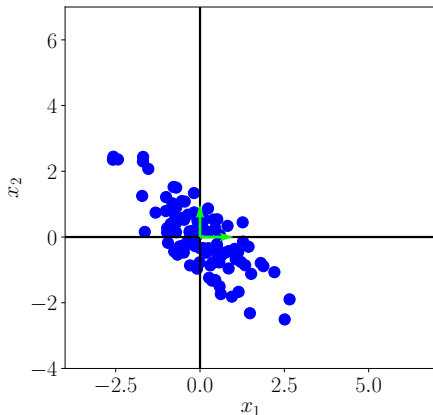
# PCA Algorithm


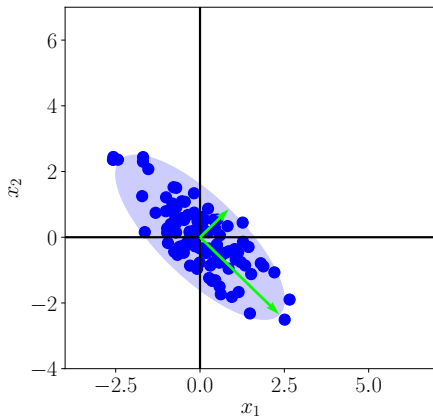
▸ Dataset

# PCA Algorithm: Step 1



‣ Mean subtraction

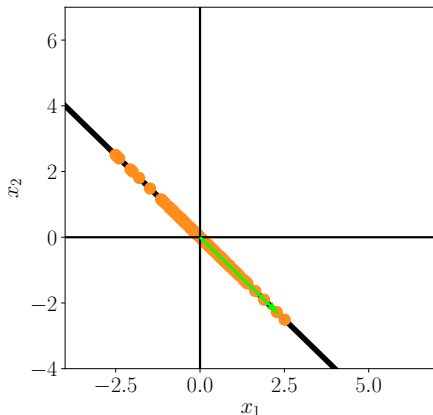# PCA Algorithm: Step 2



▸ Standardization (variance 1 in each direction)
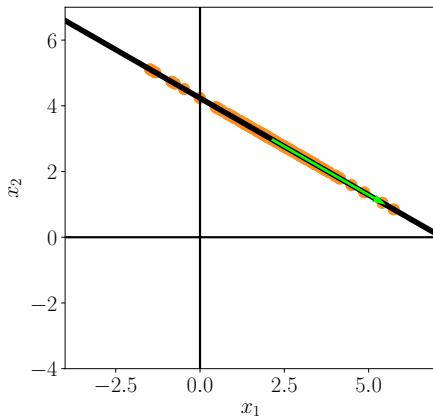
# PCA Algorithm: Step 3



▸ Eigendecomposition of the data covariance matrix

# PCA Algorithm: Step 4



▸ Orthogonal projection onto the principal subspace

# PCA Algorithm: Step 5



- Moving back to the original data space

# Overview

# PCA for High-Dimensional Data

- Fewer data points than dimensions, i.e., $N < D$.
- At least $D - N + 1$ eigenvalues 0.
- Computation time for computing eigenvalues of data covariance matrix $S$: $\mathcal{O}(D^3)$
- Rephrase PCA

# Reformulating PCA

‣ Define $X$ to be the $D \times N$-dimensional centered data matrix, whose $n$th row is $(x_n - \mathbb{E}[x])^\top$ ▶ Mean normalization

# Reformulating PCA

- Define $X$ to be the $D \times N$-dimensional centered data matrix, whose $n$th row is $(x_n - \mathbb{E}[x])^\top$ ▶▶ Mean normalization
- Corresponding covariance: $S = \frac{1}{N} X X^\top$

# Reformulating PCA

- Define $X$ to be the $D \times N$-dimensional centered data matrix, whose $n$th row is $(x_n - \mathbb{E}[x])^\top$ ▶▶ Mean normalization
- Corresponding covariance: $S = \frac{1}{N}XX^\top$
- Corresponding eigenvector equation:

$$Sb_i = \lambda_i b_i \iff \frac{1}{N}XX^\top b_i = \lambda_i b_i$$

# Reformulating PCA

- Define $X$ to be the $D \times N$-dimensional centered data matrix, whose $n$th row is $(x_n - \mathbb{E}[x])^\top$ ▶▶ Mean normalization
- Corresponding covariance: $S = \frac{1}{N}XX^\top$
- Corresponding eigenvector equation:

$$Sb_i = \lambda_i b_i \iff \frac{1}{N}XX^\top b_i = \lambda_i b_i$$

- Transformation (left-multiply by $X^\top$):

$$\frac{1}{N}XX^\top b_i = \lambda_i b_i \iff \frac{1}{N}X^\top X \underbrace{X^\top b_i}_{=:v_i} = \lambda_i \underbrace{X^\top b_i}_{=:v_i}$$

▶▶ $v_i$ is an eigenvector of the $N \times N$-matrix $\frac{1}{N}X^\top X$, which has the same non-zero eigenvalues as the original covariance matrix.

▶▶ Get eigenvalues in $\mathcal{O}(N^3)$ instead of $\mathcal{O}(D^3)$.

# Recovering the Original Eigenvectors

▸ The new eigenvalue/eigenvector equation is

$$\frac{1}{N}\boldsymbol{X}^{\top}\boldsymbol{X}\boldsymbol{v}_i = \lambda_i\boldsymbol{v}_i$$

where $\boldsymbol{v}_i = \boldsymbol{X}^{\top}\boldsymbol{b}_i$

# Recovering the Original Eigenvectors

▸ The new eigenvalue/eigenvector equation is

$$\frac{1}{N}\boldsymbol{X}^\top \boldsymbol{X}\boldsymbol{v}_i = \lambda_i \boldsymbol{v}_i$$

where $\boldsymbol{v}_i = \boldsymbol{X}^\top \boldsymbol{b}_i$

▸ We want to recover the original eigenvectors $\boldsymbol{b}_i$ of the data covariance matrix $\boldsymbol{S} = \frac{1}{N}\boldsymbol{X}\boldsymbol{X}^\top$

# Recovering the Original Eigenvectors

‣ The new eigenvalue/eigenvector equation is

$$\frac{1}{N}\boldsymbol{X}^\top \boldsymbol{X} \boldsymbol{v}_i = \lambda_i \boldsymbol{v}_i$$

where $\boldsymbol{v}_i = \boldsymbol{X}^\top \boldsymbol{b}_i$

‣ We want to recover the original eigenvectors $\boldsymbol{b}_i$ of the data covariance matrix $\boldsymbol{S} = \frac{1}{N}\boldsymbol{X}\boldsymbol{X}^\top$

‣ Left-multiply eigenvector equation by $\boldsymbol{X}$ yields

$$\underbrace{\frac{1}{N}\boldsymbol{X}\boldsymbol{X}^\top}_{=\boldsymbol{S}} \boldsymbol{X}\boldsymbol{v}_i = \lambda_i \boldsymbol{X}\boldsymbol{v}_i$$

and we recover $\boldsymbol{X}\boldsymbol{v}_i$ as an eigenvector of $\boldsymbol{S}$ associated with eigenvalue $\lambda_i$

# Recovering the Original Eigenvectors

▸ The new eigenvalue/eigenvector equation is

$$\frac{1}{N}X^\top X v_i = \lambda_i v_i$$

where $v_i = X^\top b_i$

▸ We want to recover the original eigenvectors $b_i$ of the data covariance matrix $S = \frac{1}{N}XX^\top$

▸ Left-multiply eigenvector equation by $X$ yields

$$\underbrace{\frac{1}{N}XX^\top}_{=S} X v_i = \lambda_i X v_i$$

and we recover $X v_i$ as an eigenvector of $S$ associated with eigenvalue $\lambda_i$

▸ Make sure to normalize $X v_i$ so that $\|X v_i\| = 1$

# Overview

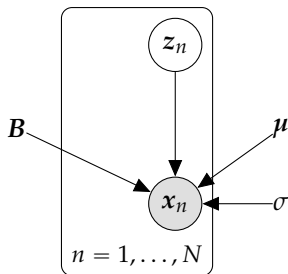# Latent Variable Perspective



▸ Model:

$$x = Bz + \mu + \epsilon, \quad \epsilon \sim \mathcal{N}\left(0,\, \sigma^2 I\right)$$

# Latent Variable Perspective



▸ Model:

$$x = Bz + \mu + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

▸ Generative process:

$$z \sim \mathcal{N}(0, I)$$
$$x|z \sim \mathcal{N}(x \mid Bz + \mu, \sigma^2 I)$$

# Why is this useful?

▸ "Standard" PCA as a special case,

# Why is this useful?

- ▸ "Standard" PCA as a special case,
- ▸ Comes with a likelihood function, and we can explicitly deal with noisy observations

# Why is this useful?

- "Standard" PCA as a special case,
- Comes with a likelihood function, and we can explicitly deal with noisy observations
- Allow for Bayesian model comparison via the marginal likelihood

# Why is this useful?

- ▸ "Standard" PCA as a special case,
- ▸ Comes with a likelihood function, and we can explicitly deal with noisy observations
- ▸ Allow for Bayesian model comparison via the marginal likelihood
- ▸ PCA as a generative model, which allows us to simulate new data

# Why is this useful?

- "Standard" PCA as a special case,
- Comes with a likelihood function, and we can explicitly deal with noisy observations
- Allow for Bayesian model comparison via the marginal likelihood
- PCA as a generative model, which allows us to simulate new data
- Straightforward connections to related algorithms and models (e.g., ICA)

# Why is this useful?

- "Standard" PCA as a special case,
- Comes with a likelihood function, and we can explicitly deal with noisy observations
- Allow for Bayesian model comparison via the marginal likelihood
- PCA as a generative model, which allows us to simulate new data
- Straightforward connections to related algorithms and models (e.g., ICA)
- Deal with data dimensions that are missing at random by applying Bayes' theorem

# Likelihood

▸ Model:

$$x = Bz + \mu + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

# Likelihood

▸ Model:
$$x = Bz + \mu + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

▸ PPCA Likelihood (integrate out the latent variables):
$$p(x|B, \mu, \sigma^2) = \int p(x|z, \mu, \sigma^2) p(z) dz$$
$$= \int \mathcal{N}(x \,|\, Bz + \mu, \sigma^2 I) \mathcal{N}(z \,|\, 0, I) dz$$

# Likelihood

‣ Model:
$$x = Bz + \mu + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

‣ PPCA Likelihood (integrate out the latent variables):
$$p(x|B, \mu, \sigma^2) = \int p(x|z, \mu, \sigma^2) p(z) dz$$
$$= \int \mathcal{N}(x \mid Bz + \mu, \sigma^2 I) \mathcal{N}(z \mid 0, I) dz$$

▶▶ Is Gaussian with mean and covariance

# Likelihood

▸ Model:

$$x = Bz + \mu + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

▸ PPCA Likelihood (integrate out the latent variables):

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

▶▶ Is Gaussian with mean and covariance

$$
\begin{aligned}
\mathbb{E}[x] &= \mathbb{E}_z[Bz + \mu + \epsilon] = \mu \\
\mathbb{V}[x] &= \mathbb{V}_z[Bz + \mu + \epsilon] = BB^\top + \sigma^2 I
\end{aligned}
$$

# Joint Distribution and Posterior

‣ Joint distribution of observed and latent variables

$$p(\boldsymbol{x}, \boldsymbol{z} | \boldsymbol{B}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N} \left( \begin{bmatrix} \boldsymbol{x} \\ \boldsymbol{z} \end{bmatrix} \middle| \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{B}\boldsymbol{B}^\top + \sigma^2 \boldsymbol{I} & \boldsymbol{B} \\ \boldsymbol{B}^\top & \boldsymbol{I} \end{bmatrix} \right)$$

‣ Posterior via Gaussian conditioning:

$$p(\boldsymbol{z} | \boldsymbol{x}, \boldsymbol{B}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\boldsymbol{z} \,|\, \boldsymbol{m}, \boldsymbol{C})$$
$$\boldsymbol{m} = \boldsymbol{B}^\top (\boldsymbol{B}\boldsymbol{B}^\top + \sigma^2 \boldsymbol{I})^{-1} (\boldsymbol{x} - \boldsymbol{\mu})$$
$$\boldsymbol{C} = \boldsymbol{I} - \boldsymbol{B}^\top (\boldsymbol{B}\boldsymbol{B}^\top + \sigma^2 \boldsymbol{I})^{-1} \boldsymbol{B}$$

▶▶ For a new observation $\boldsymbol{x}_*$ compute the posterior on $p(\boldsymbol{z}_* | \boldsymbol{x}_*, \boldsymbol{X})$ and examine it (e.g., variance).

‣ Generate new (plausible) data from this posterior

# Maximum Likelihood Estimation

▸ In PPCA, we can determine the parameters $\boldsymbol{\mu}, \boldsymbol{B}, \sigma^2$ via maximum likelihood estimation. PPCA Likelihood: $p(\boldsymbol{X}|\boldsymbol{\mu}, \boldsymbol{B}, \sigma^2)$

▸ Result (e.g., Tipping & Bishop (1999)):

$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N}\sum_{n=1}^{N} \boldsymbol{x}_n \quad \blacktriangleright\!\blacktriangleright \text{Sample mean}$$

$$\boldsymbol{B}_{\text{ML}} = \boldsymbol{T}(\boldsymbol{\Lambda} - \sigma^2 \boldsymbol{I})^{\frac{1}{2}} \boldsymbol{R}$$

$$\sigma_{\text{ML}}^2 = \frac{1}{D-M}\sum_{j=M+1}^{D} \lambda_j \quad \blacktriangleright\!\blacktriangleright \text{Average variance in orth. complement}$$

▸ For $\sigma \to 0$ the maximum likelihood solution gives the same result as PCA (see mml-book.com)

# Overview

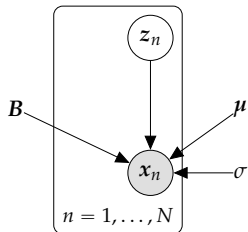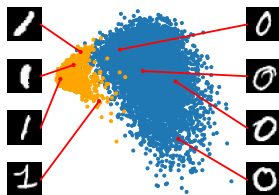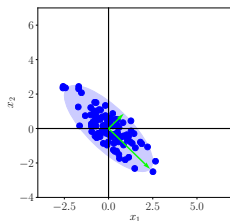# Related Models

- Factor analysis:
  Axis-aligned noise
  (instead of isotropic)

- Independent
  component analysis:
  Non-Gaussian prior
  $p(z) = \prod_m p_m(z_m)$

- Kernel PCA

- Bayesian PCA: Priors
  on parameters $B, \mu, \sigma^2$
  ▶ Approximate
  inference



- Gaussian process latent variable
  model (GP-LVM): Replace linear
  mapping in Bayesian PCA with
  Gaussian process. Point estimate of $z$

- Bayesian GP-LVM maintains a
  distribution on $z$ ▶ Approximate
  inference

# Summary



- PCA: Algorithm for linear dimensionality reduction
- Orthogonal projection of data onto a lower-dimensional subspace

    - Maximizes the variance of the projection
    - Minimizes the average squared projection/reconstruction error

- High-dimensional data
- Probabilistic PCA