



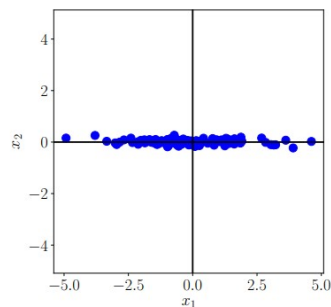
# **Chapter 10: Dimensionality Reduction with Principal Component Analysis**

Mathematics for Machine Learning

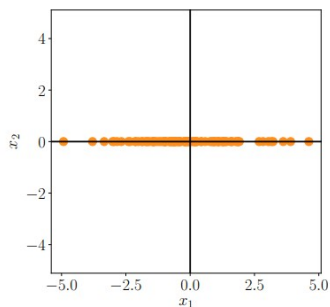
San Diego Machine Learning  
Liam Barstad

# Principal Component Analysis (PCA)

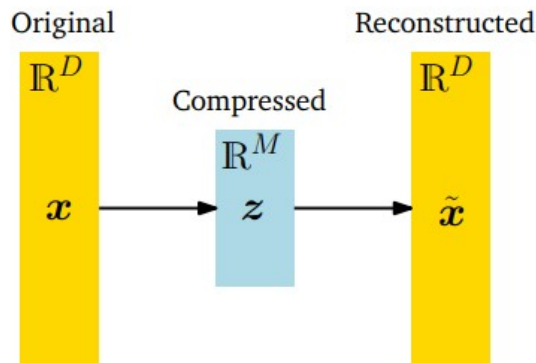
Linear method for **dimensionality reduction**



(a) Dataset with  $x_1$  and  $x_2$  coordinates.



(b) Compressed dataset where only the  $x_1$  coordinate is relevant.



$$z_n = B^\top x_n \in \mathbb{R}^M$$

$B$  = orthonormal basis (ONB) for the projection of  $x$  to  $z$ , defining the **principal subspace**  $U$

How do we minimize compression loss?

- Maximize variance in  $z$
- Minimize reconstruction loss

# Maximizing Variance Perspective

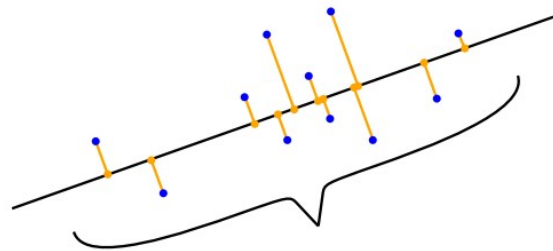
Finding the values for  $B$  (the principal subspace ONB) and  $z$  (the coordinates of the projection) that maximize variance retains the most information about  $X$

Variance for first  $z$  coordinate:  $V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2$

Substitute for covariance matrix:

$$V_1 = \frac{1}{N} \sum_{n=1}^N (b_1^\top x_n)^2 = \frac{1}{N} \sum_{n=1}^N b_1^\top x_n x_n^\top b_1$$

$$S = \frac{1}{N} \sum_{n=1}^N x_n x_n^\top \quad b_1^\top \left( \frac{1}{N} \sum_{n=1}^N x_n x_n^\top \right) b_1 = b_1^\top S b_1,$$



Re-write as optimization problem:

$$\max_{b_1} b_1^\top S b_1$$

$$\text{subject to } \|b_1\|^2 = 1.$$

$$\mathcal{L}(b_1, \lambda) = b_1^\top S b_1 + \lambda_1 (1 - b_1^\top b_1)$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = 2b_1^\top S - 2\lambda_1 b_1^\top, \quad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - b_1^\top b_1,$$

$$S b_1 = \lambda_1 b_1,$$

$$b_1^\top b_1 = 1.$$

$$\begin{array}{c} \textcolor{green}{A} \\ \text{---} \\ \text{n} \times \text{n} \\ \text{Matrix} \end{array} \begin{array}{c} \textcolor{red}{x} \\ \text{---} \\ \text{Eigenvector} \end{array} = \begin{array}{c} \textcolor{blue}{\lambda} \\ \text{---} \\ \text{Eigenvalue} \end{array} \begin{array}{c} \textcolor{red}{x} \\ \text{---} \\ \text{Eigenvector} \end{array} \quad \infty$$

# Projection Perspective

Can also be thought of as minimizing reconstruction error, distance from points to projection

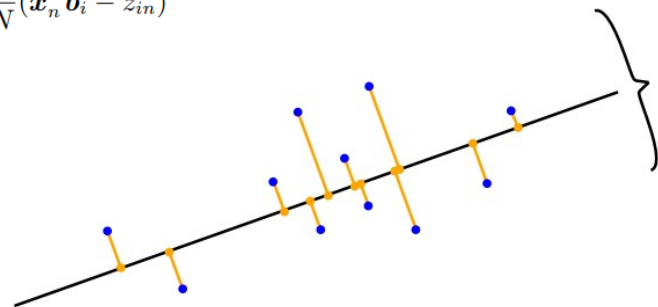
$$J_M := \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2,$$

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

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

$$\frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}} \stackrel{(10.28)}{=} \frac{\partial}{\partial z_{in}} \left( \sum_{m=1}^M z_{mn} \mathbf{b}_m \right) = \mathbf{b}_i$$

$$\stackrel{\text{ONB}}{=} -\frac{2}{N} (\mathbf{x}_n^\top \mathbf{b}_i - z_{in} \mathbf{b}_i^\top \mathbf{b}_i) = -\frac{2}{N} (\mathbf{x}_n^\top \mathbf{b}_i - z_{in})$$



Coordinates :  $z_{in} = \mathbf{x}_n^\top \mathbf{b}_i = \mathbf{b}_i^\top \mathbf{x}_n$

Optimal projection is orthonormal:

$$\tilde{\mathbf{x}}_n = \sum_{m=1}^M z_{mn} \mathbf{b}_m \stackrel{(10.32)}{=} \sum_{m=1}^M (\mathbf{x}_n^\top \mathbf{b}_m) \mathbf{b}_m.$$

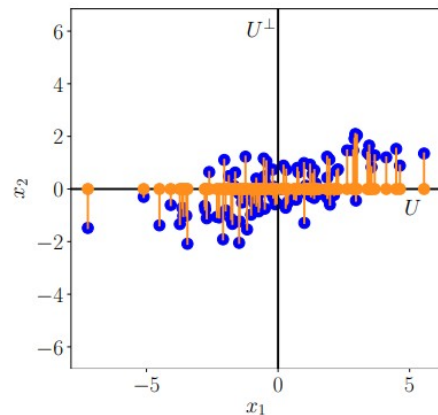
$$\mathbf{x}_n = \sum_{d=1}^D z_{dn} \mathbf{b}_d \stackrel{(10.32)}{=} \sum_{d=1}^D (\mathbf{x}_n^\top \mathbf{b}_d) \mathbf{b}_d = \left( \sum_{d=1}^D \mathbf{b}_d \mathbf{b}_d^\top \right) \mathbf{x}_n$$

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

$$\tilde{\mathbf{x}} = \mathbf{B} \underbrace{(\mathbf{B}^\top \mathbf{B})^{-1}}_{=\mathbf{I}} \mathbf{B}^\top \mathbf{x} = \mathbf{B} \mathbf{B}^\top \mathbf{x},$$

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

$$= \sum_{j=M+1}^D (\mathbf{x}_n^\top \mathbf{b}_j) \mathbf{b}_j.$$



# Projection Perspective - Basis Vectors

We get projection matrix:  $\sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^\top = \mathbf{B} \mathbf{B}^\top.$

In order to minimize, need to find the best rank- $M$  approximation of  $\mathbf{B} \mathbf{B}^\top = \mathbf{I}$

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 &= \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{B} \mathbf{B}^\top \mathbf{x}_n\|^2 \\ &= \frac{1}{N} \sum_{n=1}^N \|( \mathbf{I} - \mathbf{B} \mathbf{B}^\top ) \mathbf{x}_n\|^2. \end{aligned}$$

$$J_M = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 \stackrel{(10.38b)}{=} \frac{1}{N} \sum_{n=1}^N \left\| \sum_{j=M+1}^D (\mathbf{b}_j^\top \mathbf{x}_n) \mathbf{b}_j \right\|^2$$

$$\begin{aligned} J_M &= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (\mathbf{b}_j^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{x}_n \mathbf{b}_j^\top \mathbf{x}_n \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_j, \end{aligned}$$

$$\begin{aligned} J_M &= \sum_{j=M+1}^D \mathbf{b}_j^\top \underbrace{\left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)}_{=: \mathbf{S}} \mathbf{b}_j = \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{S} \mathbf{b}_j \quad (10.43a) \\ &= \sum_{j=M+1}^D \text{tr}(\mathbf{b}_j^\top \mathbf{S} \mathbf{b}_j) = \sum_{j=M+1}^D \text{tr}(\mathbf{S} \mathbf{b}_j \mathbf{b}_j^\top) = \text{tr} \left( \underbrace{\left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top \right)}_{\text{projection matrix}} \mathbf{S} \right), \end{aligned}$$

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

The distance between the data's subspace and the projection's subspace is proportional to the covariance matrix  $\mathbf{S}$ , because the optimal projection is orthogonal

From here, the derivations can continue from the maximum variance calculations

Therefore, minimizing the projection is equivalent to maximizing the variance

# Summary So Far

- **Principal components** - eigenvectors of the correlation matrix  $S$
- PCA uses the first  $M$  principal components to generate the basis for a portion of the underlying variance of the data
- A new subspace, the principal subspace, is created from the principal components, and data points are projected from the original set using the linear transformation  $z = B^T x$
- This projection both minimizes the euclidean distance between the data and its projection and maximizes the variance of its encoding

# Methods of Computing PCA

- Can perform eigendecomposition since  $S$  is square
- Can use SVD – Columns of  $U$  are eigenvectors of  $XX^T$ , or  $S$

$$\underbrace{\mathbf{X}}_{D \times N} = \underbrace{\mathbf{U}}_{D \times D} \underbrace{\mathbf{\Sigma}}_{D \times N} \underbrace{\mathbf{V}^T}_{N \times N},$$

This relationship between the eigenvalues of  $S$  and the singular values of  $\mathbf{X}$  provides the connection between the maximum variance view (Section 10.2) and the singular value decomposition.

- Using low rank approximations  $\tilde{\mathbf{X}}_M := \operatorname{argmin}_{\operatorname{rk}(\mathbf{A}) \leq M} \|\mathbf{X} - \mathbf{A}\|_2 \in \mathbb{R}^{D \times N}$
- Power Iteration (for first eigenvector)  $\mathbf{x}_{k+1} = \frac{S\mathbf{x}_k}{\|S\mathbf{x}_k\|}, \quad k = 0, 1, \dots$
- For high dimensional data (e.g. images) can turn  $D \times D$  into  $N \times N$  matrix

$$S\mathbf{b}_m = \lambda_m \mathbf{b}_m, \quad m = 1, \dots, M,$$

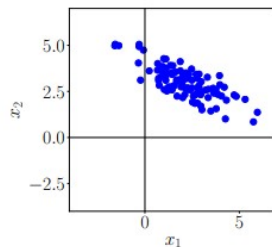
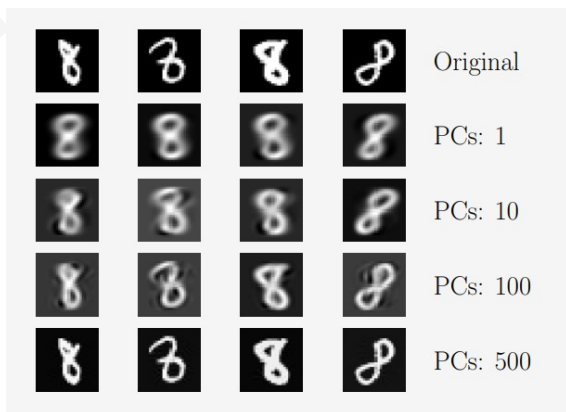
$$\frac{1}{N} \underbrace{\mathbf{X}^T \mathbf{X}}_{N \times N} \underbrace{\mathbf{X}^T \mathbf{b}_m}_{=: \mathbf{c}_m} = \lambda_m \mathbf{X}^T \mathbf{b}_m \iff \frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{c}_m,$$

- Get eigenvectors of  $X^T X$ , then recover  $X\mathbf{c}_m$  as eigenvector of  $S$

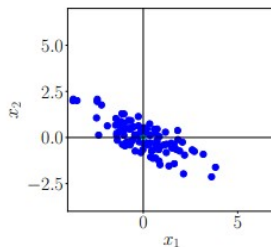
$$\underbrace{\frac{1}{N} \mathbf{X} \mathbf{X}^T}_S \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{X} \mathbf{c}_m$$

# Steps to Compute PCA

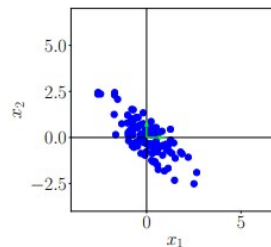
- Subtract mean from all data points (mean = 0)
- **Standardize** (divide data points by stddev)
- Get eigenvectors of covariance matrix
- Project data onto principal subspace
- Multiply by original stddev and add back mean



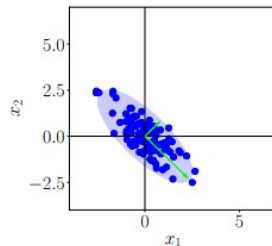
(a) Original dataset.



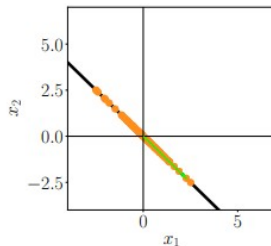
(b) Step 1: Centering by subtracting the mean from each data point.



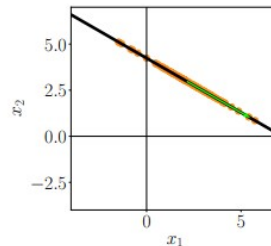
(c) Step 2: Dividing by the standard deviation to make the data unit free. Data has variance 1 along each axis.



(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 PCA (PPCA)

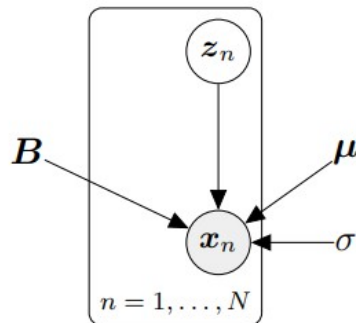
Can also consider  $z$  to be a latent variable

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

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

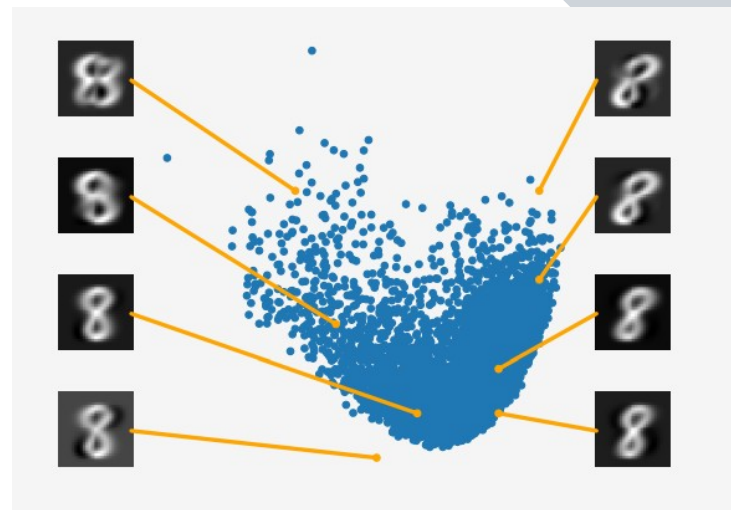
$$\mathbf{z}_n \sim \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$$

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



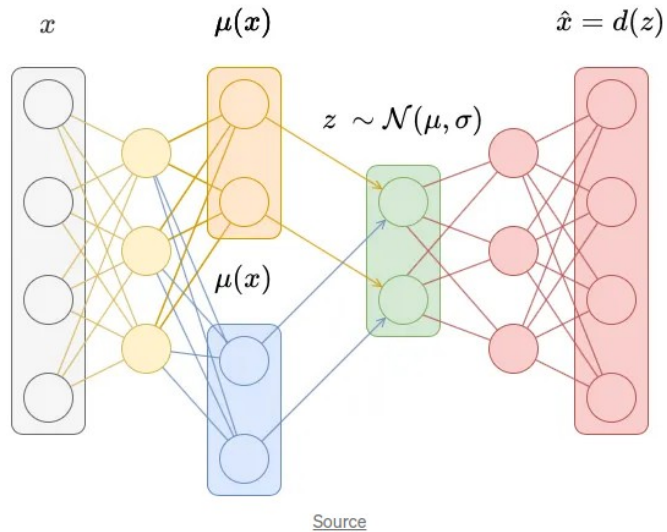
Can use likelihood and posterior for MAP and MLE

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



# Example - Variational Autoencoder

Trains statistics for normally distributed latent variable  $z$



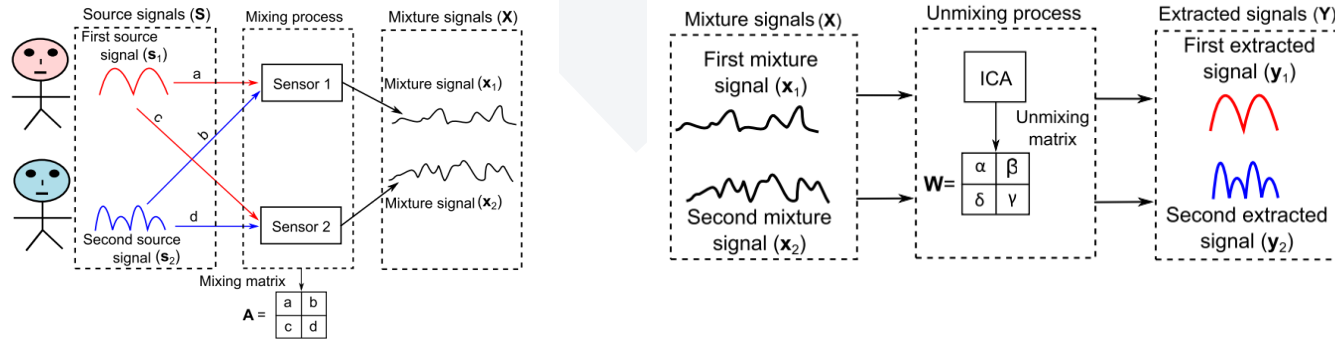
- Optimizes objective function (euclidean distance, but also implicitly variance)
- Lower dimensional representation

<https://medium.com/@weidagang/demystifying-neural-networks-variational-autoencoders-6a44e75d0271>

# Further Reading - Independent Component Analysis

PCA models lower dimensional subspace, ICA models underlying signals

PCA optimizes covariance, ICA optimizes higher-order metrics like kurtosis



ICA extracts I.I.D non-Gaussian signals and measure “Gaussianity”

## Steps

- Subtract mean from data points
- “decorrelate” – project into PCA principal subspace
- Scale/normalize
- Find the “unmixing” matrix  $W$ , where  $x = ASW$ , by maximizing “non-gaussianity”, e.g. kurtosis, negentropy
- Projecting  $x$  into  $W$  gives signals

<https://www.emerald.com/insight/content/doi/10.1016/j.aci.2018.08.006/full/pdf?title=independent-component-analysis-an-introduction>