

Mathematics for Machine Learning

Prof Willie Brink

Applied Mathematics, Stellenbosch University

Lecture 9: Dimensionality Reduction with PCA

Contents of the module

Chapter 02: Linear Algebra

Chapter 03: Analytic Geometry

Chapter 04: Matrix Decompositions

Chapter 05: Vector Calculus

Chapter 06: Probability and Distributions

Chapter 07: Continuous Optimisation

Chapter 08: When Models Meet Data

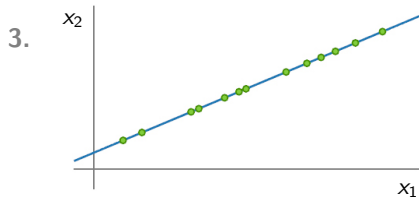
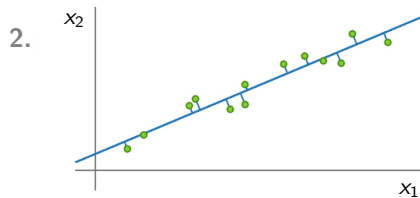
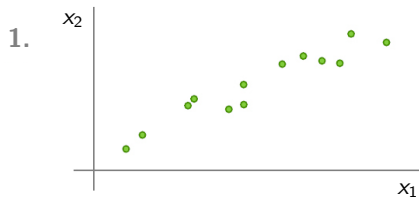
Chapter 09: Linear Regression

Chapter 10: Dimensionality Reduction with Principal Component Analysis

Chapter 11: Density Estimation with Gaussian Mixture Models

Chapter 12: Classification with Support Vector Machines

Introduction



10.1 Problem setting

High-dimensional data is often overcomplete (many redundant dimensions), and may occupy a much lower-dimensional subspace.

Consider data points $\mathbf{x}, \dots, \mathbf{x}_N$ in \mathbb{R}^D , with a mean of $\mathbf{0}$.

PCA: find projections $\tilde{\mathbf{x}}_n$ of data points \mathbf{x}_n , that are similar to original data but have a significantly lower intrinsic dimensionality.

Let $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$ be a projection matrix with orthonormal columns, where $M \ll D$. Our task will be to find this matrix \mathbf{B} for a given dataset.

Encoding \mathbf{x}_n to a low-dimensional representation: $\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M$

Decoding \mathbf{z}_n in order to reconstruct \mathbf{x}_n : $\tilde{\mathbf{x}}_n = \mathbf{B} \mathbf{z}_n = \mathbf{B} \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^D$

10.2 Maximum variance perspective

Find a \mathbf{B} that retains as much information as possible, i.e. captures the most variance.

Let's start by finding a single vector $\mathbf{b}_1 \in \mathbb{R}^D$ that maximises the variance of the first coordinate z_1 of the encodings (the first principal component).

$$V_1 = \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1,n}^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{b}_1^T \mathbf{x}_n)^2 = \mathbf{b}_1^T \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_1 = \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1$$

Solve the constrained optimisation problem: $\max_{\mathbf{b}_1} \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1$ subject to $\|\mathbf{b}_1\|^2 = 1$

Introducing a Lagrange multiplier λ_1 and setting derivatives w.r.t. \mathbf{b}_1 and λ_1 to 0, give

$$\mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1, \quad \mathbf{b}_1^T \mathbf{b}_1 = 1 \quad \text{and note then that } V_1 = \lambda_1 \mathbf{b}_1^T \mathbf{b}_1 = \lambda_1$$

Therefore we choose \mathbf{b}_1 as the eigenvector of \mathbf{S} associated with its largest eigenvalue.

\mathbf{b}_1 is the eigenvector of data covariance matrix \mathbf{S} associated with its largest eigenvalue.

The second principal component, \mathbf{b}_2 , will be the eigenvector of \mathbf{S} associated with the second largest eigenvalue, and so on.

The first M principal components form an ONB for an M -dimensional subspace of \mathbb{R}^D .

The maximum amount of variance that PCA can capture is $V_M = \sum_{m=1}^M \lambda_m$

The variance lost by PCA's compression is $J_M = \sum_{j=M+1}^D \lambda_j = V_D - V_M$

Note: if the data is not centered at $\mathbf{0}$, we would first subtract the data mean $\boldsymbol{\mu}$ from each \mathbf{x}_n before forming \mathbf{S} and finding \mathbf{B} .

The encoding would then be $\mathbf{z}_n = \mathbf{B}^T(\mathbf{x}_n - \boldsymbol{\mu})$, and the decoding would be $\tilde{\mathbf{x}}_n = \mathbf{B}\mathbf{z}_n + \boldsymbol{\mu}$.

10.3 Projection perspective

PCA can also be derived from the perspective of a linear encoder-decoder that minimises the average **reconstruction error**.

The aim is to find vectors $\tilde{\mathbf{x}}_n \in \mathbb{R}^D$ that lie in an M -dimensional subspace spanned by an unknown ONB $(\mathbf{b}_1, \dots, \mathbf{b}_M)$, that is as close as possible to the original data \mathbf{x}_n , i.e. that minimise the average reconstruction error:

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

It turns out that for a given ONB an **orthogonal projection** gives the optimal encoding.

It also turns out that minimising the reconstruction error is equivalent to **minimising the variance we ignore** when projecting to the subspace, leading to the same solution as before (eigenvectors of the data covariance matrix \mathbf{S}).

10.4 Eigenvector computation

Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}_{D \times N}$. We obtain the principal components as eigenvectors of the data covariance matrix \mathbf{S} , where

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

Recall that the first M cols of \mathbf{U} in the SVD of \mathbf{X} give us exactly those eigenvectors!

The eigenvalues λ_m of \mathbf{S} are related to the singular values σ_m of \mathbf{X} via: $\lambda_m = \sigma_m^2 / N$.

We would normally use the SVD of \mathbf{X} to perform PCA, for its numerical stability and computational efficiency (compared to an eigendecomposition of \mathbf{S}).

10.6 Key steps of PCA in practice

Given a dataset of points $\mathbf{x}_1, \dots, \mathbf{x}_N$ in \mathbb{R}^D .

1. Mean subtraction

Center the data at $\mathbf{0}$ by subtracting the mean $\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$ from each \mathbf{x}_n .

2. Standardisation

Divide each data point by the standard deviation σ_d for every dimension $d = 1, \dots, D$.
Now the data has variance 1 along each axis.

3. Determining the principal components

Concatenate the centered, standardised data vectors as columns of \mathbf{X} , and let \mathbf{B} be the first M columns of \mathbf{U} in the SVD of \mathbf{X} .

4. Projection

Any point $\mathbf{x} \in \mathbb{R}^D$ (from the same data generating process as the given dataset) can be encoded as a lower-dimensional vector: $\mathbf{z} = \mathbf{B}^T \mathbf{x}_*$, where \mathbf{x}_* has components

$$x_*^{(d)} = \frac{x^{(d)} - \mu_d}{\sigma_d}$$

The vector \mathbf{z} is an M -dimensional representation of the D -dimensional vector \mathbf{x} .

5. Reconstruction

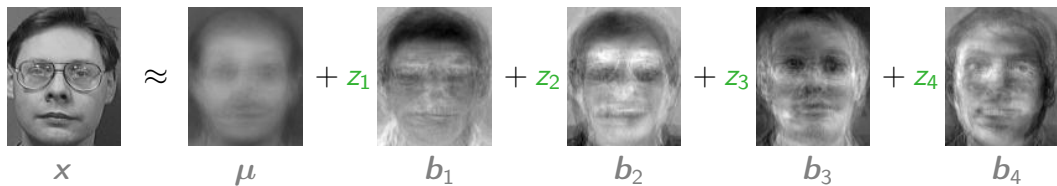
A representation \mathbf{z} is transformed back to D -dimensional space by $\tilde{\mathbf{x}}_* = \mathbf{B}\mathbf{z}$, and then de-standardising:

$$\tilde{x}^{(d)} = \tilde{x}_*^{(d)} \sigma_d + \mu_d$$

The vector $\tilde{\mathbf{x}}$ might be an approximation of the original \mathbf{x} from step 4.

PCA on an image dataset

Dataset:



$$z = [z_1 \ z_2 \ z_3 \ z_4]^T$$