

Machine Learning

Deniz Gündüz and Krystian Mikołajczyk

Department of Electrical and Electronic Engineering
Imperial College London

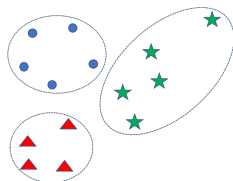
- Unsupervised learning
- Linear autoencoder
- Principal component analysis (PCA)

Machine Learning

- **Supervised learning:** Given data samples with labels (\mathbf{x}, y) , we want to learn a function $y = f(\mathbf{x})$ to predict labels of new samples
 - ▶ Predictive analytics
 - ▶ Classification: y is discrete
 - ▶ Regression: y is continuous
- **Unsupervised learning:** We are given only samples of data $\mathbf{x}_1, \dots, \mathbf{x}_n$, we want to compute a function $y = f(\mathbf{x})$ that provides a simpler representation
 - ▶ Descriptive analytics
 - ▶ y is discrete: clustering
 - ▶ y is continuous: dimension reduction, autoencoders

Unsupervised Learning

- Data without labels: $\mathbf{x}_1, \dots, \mathbf{x}_n$
- **Dimension reduction**
 - ▶ Preprocessing step for supervised learning
 - ▶ Visualisation



- **Clustering**
 - ▶ Data compression
 - ▶ Group data samples that share “similarity”

Spanning Set

- Data points: $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- Assume *mean-centered* data points: We can subtract the mean in each dimension if they are not - reversible operation
- Consider set of basis vectors $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_k\}$, $\mathbf{c}_j \in \mathbb{R}^d$, $\|\mathbf{c}_j\|^2 = 1$, $\forall j$
- \mathcal{C} is called a **spanning set** if each point \mathbf{x}_i can be written as some linear combination of basis vectors in \mathcal{C} , i.e.,

$$\mathbf{x}_i = \sum_{j=1}^k w_{i,j} \mathbf{c}_j$$

for some $w_{i,1}, \dots, w_{i,k}$.

- We can write more compactly as $\mathbf{C} \mathbf{w}_i = \mathbf{x}_i$, where

$$\mathbf{C} = [\mathbf{c}_1 \cdots \mathbf{c}_k], \quad \text{and} \quad \mathbf{w}_i = [w_{i,1}, \dots, w_{i,k}]^T$$

- \mathbf{w}_i is called **encoding** of \mathbf{x}_i over \mathcal{C} , while $\mathbf{C} \mathbf{w}_i$ is called **decoding**

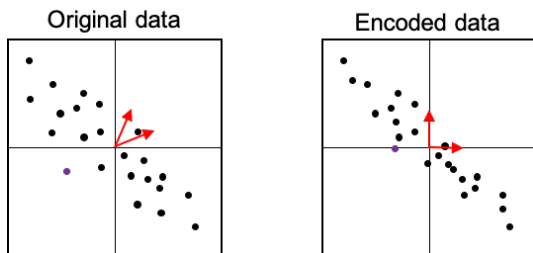
Spanning Set

If

- ① $k = d$, and
- ② vectors in \mathcal{C} are linearly independent,

then \mathcal{C} spans \mathbb{R}^d , i.e., any $\mathbf{x} \in \mathbb{R}^d$ can be written as $\mathbf{C}\mathbf{w} = \mathbf{x}$ for some \mathbf{w} .

Hence, $\mathbf{x}_i \rightarrow \mathbf{w}_i$ is simply a coordinate rotation



Orthonormal Basis

- Set $\{\mathbf{c}_1, \dots, \mathbf{c}_k\}$ is an orthonormal basis if

$$\mathbf{c}_i^T \mathbf{c}_j = \begin{cases} 1 & \text{if } i = j; \\ 0 & \text{if } i \neq j. \end{cases}$$

- Equivalently: $C^T C = I_{d \times d}$
- Encoding is simply: $\mathbf{w}_i = C^T \mathbf{x}_i$
- **Autoencoder** formula:

$$C \underbrace{C^T \mathbf{x}_i}_{\mathbf{w}_i} = \mathbf{x}_i$$

Imperfect Representation of Data

- If $k < d$, we cannot represent every possible data point in \mathbb{R}^d . We want $C\mathbf{w}_i \approx \mathbf{x}_i$

- We try to minimize the average error after decoding, i.e.,

$$\|C\mathbf{w}_i - \mathbf{x}_i\|^2$$

- If C is orthonormal, taking the derivative w.r.t. \mathbf{w}_i , we find

$$w_{i,j} = \mathbf{c}_j^T \mathbf{x}_i$$

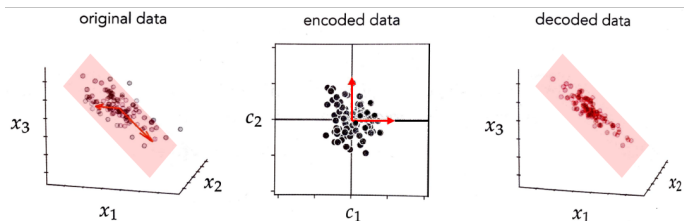


Learning a Spanning Set

- Typically we are not given a spanning set, but need to learn from data
- Given dataset $\mathbf{x}_1, \dots, \mathbf{x}_n$, we want to find the spanning set $\{\mathbf{c}_1, \dots, \mathbf{c}_k\}$ that minimizes the average error

$$\min_{C, \mathbf{w}_1, \dots, \mathbf{w}_n} \frac{1}{n} \sum_{i=1}^n \|C \mathbf{w}_i - \mathbf{x}_i\|^2$$

- Not a convex problem, but can be solved using gradient descent



Linear Autoencoder

- Remember: for orthonormal spanning vectors, we have $\mathbf{w}_i = C^T \mathbf{x}_i$
- Optimization problem can be written as

$$\min_{C: C^T C = I_{d \times d}} \frac{1}{n} \sum_{i=1}^n \|CC^T \mathbf{x}_i - \mathbf{x}_i\|^2$$

- No need to impose orthonormality

Optimal Linear Autoencoder

$$\min_C \frac{1}{n} \sum_{i=1}^n \|CC^T \mathbf{x}_i - \mathbf{x}_i\|^2$$

- CC^T is a symmetric $d \times d$ matrix
- Eigenvalue decomposition: $CC^T = VDV^T$, where $V \in \mathbb{R}^{d \times d}$ is an orthonormal matrix of eigenvectors (i.e., $VV^T = I_{d \times d}$), and $D \in \mathbb{R}^{d \times d}$ is a diagonal matrix with at most k nonnegative eigenvalues

Optimal Linear Autoencoder

$$\begin{aligned}\|CC^T \mathbf{x}_i - \mathbf{x}_i\|^2 &= \|VDV^T \mathbf{x}_i - \mathbf{x}_i\|^2 \\ &= \mathbf{x}_i^T VDDV^T \mathbf{x}_i - 2\mathbf{x}_i^T VDV^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{x}_i \\ &= \mathbf{x}_i^T VDDV^T \mathbf{x}_i - 2\mathbf{x}_i^T VDV^T \mathbf{x}_i + \mathbf{x}_i^T VV^T \mathbf{x}_i\end{aligned}$$

Define $\mathbf{q}_i = V^T \mathbf{x}_i$. Above can be rewritten as

$$\begin{aligned}\mathbf{q}_i^T DD \mathbf{q}_i - 2\mathbf{q}_i^T D \mathbf{q}_i + \mathbf{q}_i^T \mathbf{q}_i &= \mathbf{q}_i^T (D^2 - 2D + I_{d \times d}) \mathbf{q}_i \\ &= \mathbf{q}_i^T (D - I_{d \times d})^2 \mathbf{q}_i\end{aligned}$$

Optimization problem

$$\min_C \frac{1}{n} \sum_{i=1}^n \mathbf{q}_i^T (D - I_{d \times d})^2 \mathbf{q}_i$$

This is minimized when D has k eigenvalues of $+1$. This means that C is an orthonormal matrix.

Principal Component Analysis (PCA)

- Widely used for dimensionality reduction, lossy data compression, feature extraction, and data visualization
- Two formulations possible, leading to the same solution:
 - ▶ **Maximum variance formulation:** Orthogonal projection of data to a lower dimensional principal subspace, such that the variance of projected data is maximized (Hotelling'33)
 - ▶ **Minimum error formulation:** Linear projection that minimizes the mean squared error between data points and their projections (Pearson'01)

PCA - Maximum Variance Formulation

Consider $k = 1$. We want a single vector \mathbf{c} to project on with $\|\mathbf{c}\|^2 = 1$.

We want to maximize

$$\frac{1}{n} \sum_{i=1}^n \left(\mathbf{c}^T \mathbf{x}_i - \mathbf{c}^T \bar{\mathbf{x}} \right)^2$$

where $\bar{\mathbf{x}}$ is the mean-value of data, which we assumed to be zero. Variance reduces to

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left(\mathbf{c}^T \mathbf{x}_i \right)^2 &= \frac{1}{n} \sum_{i=1}^n \mathbf{c}^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{c} \\ &= \mathbf{c}^T \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{c} \\ &= \mathbf{c}^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c} \end{aligned}$$

where $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_n]$ is the data matrix.

PCA - Minimum Error Formulation

Assume orthonormal spanning set with k elements. We have

$$\begin{aligned}\|CC^T \mathbf{x}_i - \mathbf{x}_i\|^2 &= \mathbf{x}_i^T CC^T CC^T \mathbf{x}_i - 2\mathbf{x}_i^T CC^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{x}_i \\ &= -\mathbf{x}_i^T CC^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{x}_i \\ &= -\|C^T \mathbf{x}_i\|^2 + \|\mathbf{x}_i\|^2\end{aligned}$$

Since \mathbf{x}_i is given, we only minimize $-\|C^T \mathbf{x}_i\|^2 = -\sum_{j=1}^k (\mathbf{c}_j^T \mathbf{x}_i)^2$

Hence, the objective function is

$$\min_{\substack{\mathbf{c}_1, \dots, \mathbf{c}_k: \\ \|\mathbf{c}_i\|^2=1, \\ \mathbf{c}_i^T \mathbf{c}_j=0, i \neq j}} -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k (\mathbf{c}_j^T \mathbf{x}_i)^2$$

We can minimize one basis vector at a time!

Principal Component Analysis (PCA)

We begin with \mathbf{c}_1 . We want

$$\min_{\mathbf{c}_1: \|\mathbf{c}_1\|^2=1} -\frac{1}{n} \sum_{i=1}^n (\mathbf{c}_1^T \mathbf{x}_i)^2$$

Note that, this measure the variance of data along the direction of \mathbf{c}_1 (because we assumed zero-mean data): **Equivalent to max. variance formulation**

We have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (\mathbf{c}_1^T \mathbf{x}_i)^2 &= \frac{1}{n} \sum_{i=1}^n \mathbf{c}_1^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{c}_1 \\ &= \mathbf{c}_1^T \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{c}_1 \\ &= \mathbf{c}_1^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_1 \end{aligned}$$

Principal Component Analysis (PCA)

$$\begin{aligned} \min_{\mathbf{c}_1} \quad & -\mathbf{c}_1^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_1 \\ \text{s.t.} \quad & \mathbf{c}_1^T \mathbf{c}_1 = 1 \end{aligned}$$

We form the Lagrangian

$$\mathcal{L}(\mathbf{c}_1, \alpha) = -\mathbf{c}_1^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_1 + \alpha (\mathbf{c}_1^T \mathbf{c}_1 - 1)$$

We have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{c}_1} = -2 \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_1 + 2\alpha \mathbf{c}_1 = 0$$

We get $\left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_1 = \alpha \mathbf{c}_1$

Implies that α is an eigenvalue of $\left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right)$ with associated eigenvector \mathbf{c}_1 .

Principal Component Analysis (PCA)

$$\begin{aligned} \min_{\mathbf{c}_1} \quad & -\mathbf{c}_1^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_1 \\ \text{s.t.} \quad & \mathbf{c}_1^T \mathbf{c}_1 = 1 \end{aligned}$$

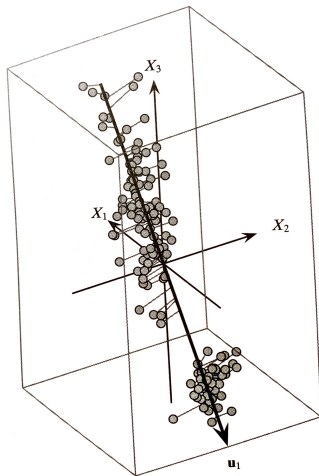
The optimal value is obtained as

$$\mathbf{c}_1^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_1 = \mathbf{c}_1^T \alpha \mathbf{c}_1 = \alpha$$

We should choose \mathbf{c}_1 as the eigenvector corresponding to the largest eigenvalue of $\left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right)$. This is called the **first principal component**.

Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ be the eigenvalues of $\frac{1}{n} \mathbf{X} \mathbf{X}^T$ with corresponding eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_k$.

Approximation with First Principal Component



Second Principal Component

$$\begin{aligned} \min_{\mathbf{c}_2} \quad & -\mathbf{c}_2^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_2 \\ \text{s.t.} \quad & \mathbf{c}_2^T \mathbf{c}_2 = 1 \text{ and } \mathbf{c}_2^T \mathbf{u}_1 = 0 \end{aligned}$$

$$\mathcal{L}(\mathbf{c}_2, \theta, \beta) = -\mathbf{c}_2^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_2 + \theta(\mathbf{c}_2^T \mathbf{c}_2 - 1) + \beta \mathbf{c}_2^T \mathbf{u}_1$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{c}_2} = -2 \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_2 + 2\theta \mathbf{c}_2 + \beta \mathbf{u}_1 = 0$$

Multiplying both sides with \mathbf{u}_1^T , we get

$$-2\mathbf{u}_1^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_2 + 2\theta \mathbf{u}_1^T \mathbf{c}_2 + \beta \mathbf{u}_1^T \mathbf{u}_1 = 0$$

$$\beta = 2\mathbf{c}_2^T \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{u}_1 = 2\mathbf{c}_2^T \lambda_1 \mathbf{u}_1 = 0$$

Second Principal Component

Plugging $\beta = 0$ into $\frac{\partial \mathcal{L}}{\partial \mathbf{c}_2} = 0$ above, we get

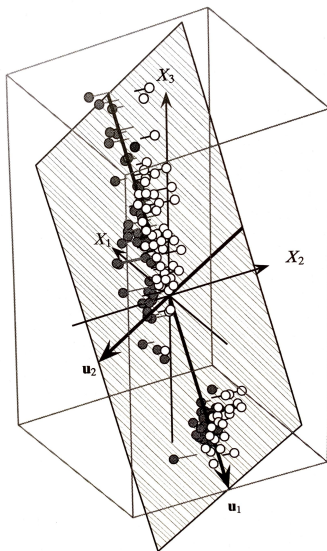
$$\left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{c}_2 = \theta \mathbf{c}_2$$

Implies that \mathbf{c}_2 is another eigenvector of $\frac{1}{n} \mathbf{X} \mathbf{X}^T$, orthogonal to \mathbf{u}_1 .

The maximum value is given by θ , its eigenvalue. Hence, the second principal component is given by $\mathbf{c}_2 = \mathbf{u}_2$.

We can continue similarly, and the j th principal component is given by the eigenvector corresponding to the j th largest eigenvalue of the mean-centered data covariance matrix $\frac{1}{n} \mathbf{X} \mathbf{X}^T$.

Approximation with Two Principal Components



Application of PCA

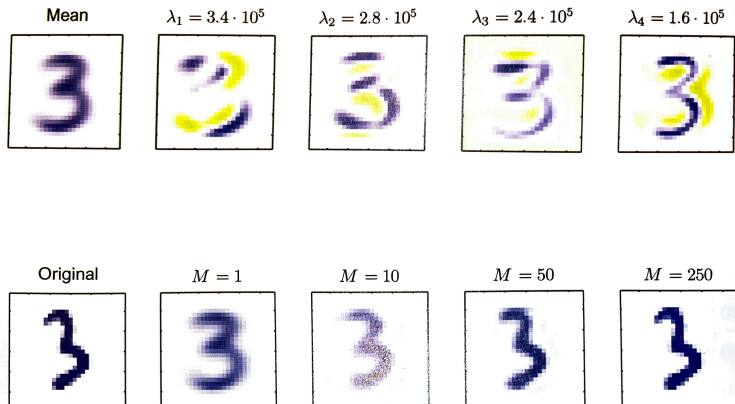


Figure from Bishop, Pattern Recognition and Machine Learning.

Application of PCA - Whitening

- Pre-process data to standardize some properties
- Let $S = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$ be the data covariance matrix ($\bar{\mathbf{x}}$ is the sample mean) with eigenvalue decomposition

$$S = ULU^T$$

where U is the orthogonal matrix with columns \mathbf{u}_i , and L is the diagonal matrix with elements λ_i .

- Transform each data points as

$$\mathbf{y}_i = L^{-1/2} U^T (\mathbf{x}_i - \bar{\mathbf{x}})$$

Application of PCA - Whitening

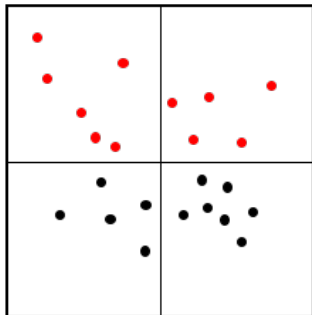
We have

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^T &= \frac{1}{n} \sum_{i=1}^n L^{-1/2} U^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T U L^{-1/2} \\&= L^{-1/2} U^T \left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \right) U L^{-1/2} \\&= L^{-1/2} U^T S U L^{-1/2} \\&= L^{-1/2} U^T \textcolor{red}{U} L U^T U L^{-1/2} \\&= L^{-1/2} L L^{-1/2} = I_{d \times d}\end{aligned}$$

Each component of the new data representation are perfectly **uncorrelated**!

Classification with PCA

Original data



Encoded data



Variance in PCA

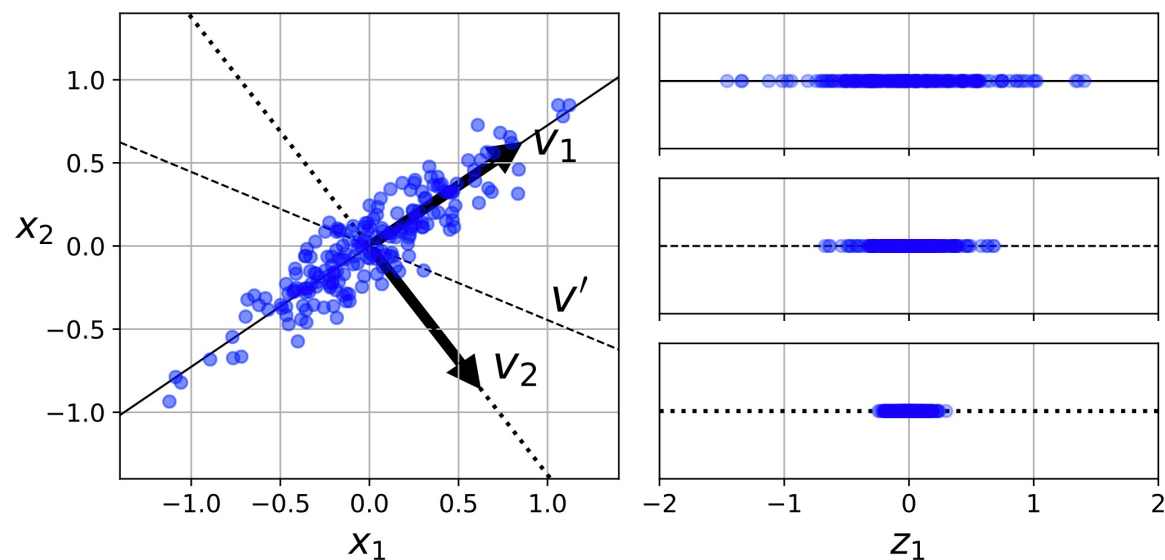
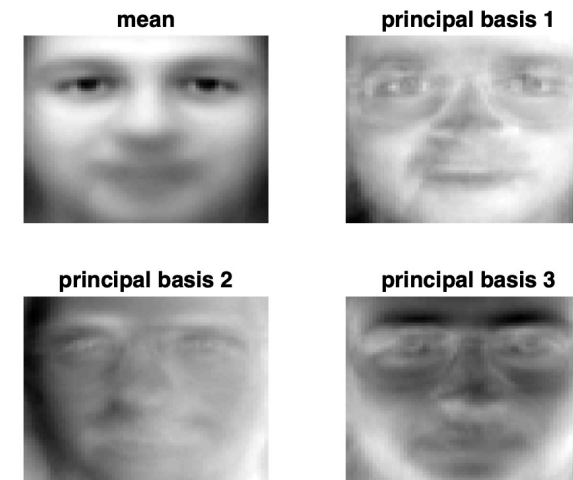


Figure 20.4: Illustration of the variance of the points projected onto different 1d vectors. v_1 is the first principal component, which maximizes the variance of the projection. v_2 is the second principal component which is direction orthogonal to v_1 . Finally v' is some other vector in between v_1 and v_2 . Adapted from Figure 8.7 of [Gér19]. Generated by code at [figures.probml.ai/book1/20.4](https://probml.ai/book1/20.4)

Example



(a)



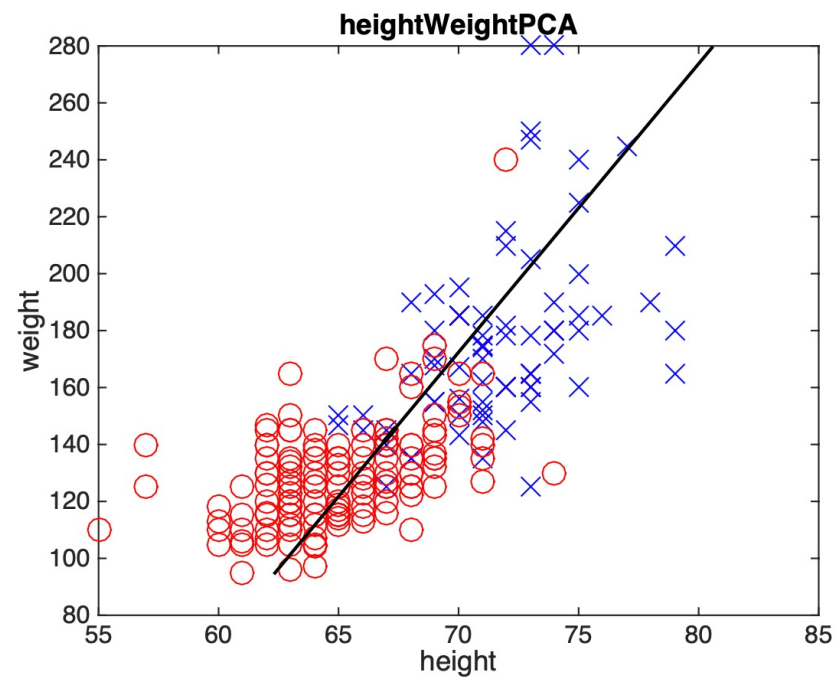
(b)

Figure 20.3: a) Some randomly chosen 64×64 pixel images from the Olivetti face database. (b) The mean and the first three PCA components represented as images. Generated by code at figures.probml.ai/book1/20.3.

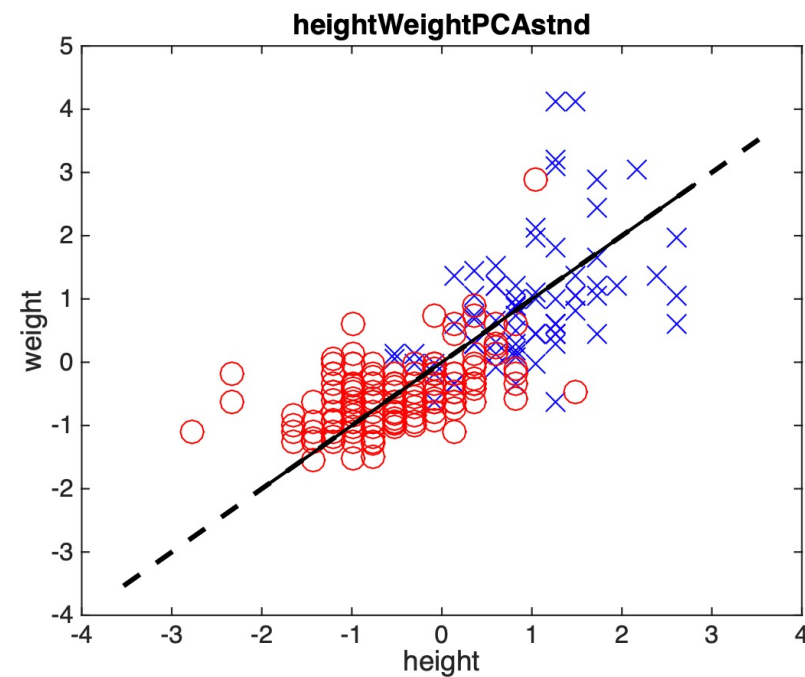
Why do we scale/normalize/standardize our data?

1. Scaling does not affect our statistical inference.
2. When we make new features out of the existing ones (e.g., by adding or multiplying two feature), then if features are not in the same scale, one of them can cancel out the other one.
3. PCA can only be interpreted as the singular value decomposition of a data matrix when the columns have first been centered by their means.
4. Simplification of analysis:
the sample covariance matrix of a matrix of values centered by their sample means is simply $X'X$. Similarly, if a univariate random variable X has been mean centered, then $\text{var}(X)=E(X^2)$

Example



(a)



(b)

Figure 20.5: Effect of standardization on PCA applied to the height/weight dataset. (Red=female, blue=male.) Left: PCA of raw data. Right: PCA of standardized data. Generated by code at figures.problml.ai/book1/20.5.