

Machine Learning and Data Mining WS21/22

# “3 Data Transformation”

Dr. Zeyd Boukhers

@ZBoukhers

Institute for Web Science and Technologies  
University of Koblenz-Landau

November 10, 2021

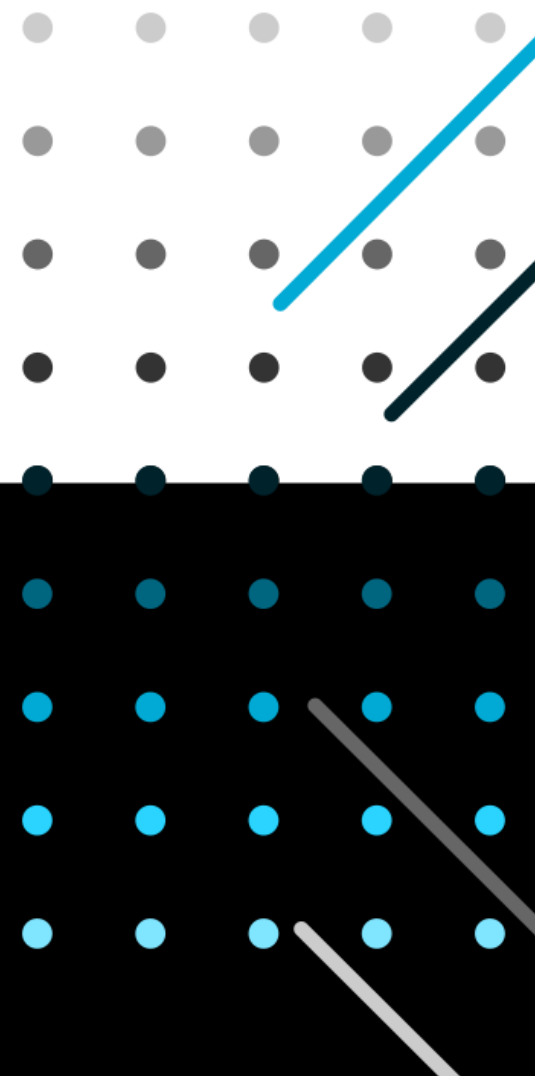


- How to define a MLDM task.
- How to design high quality features
- How to pre-process the data.
  - Remove outliers.
  - Scale features.
  - Measure the correlation between features.
  - Handle missing data.

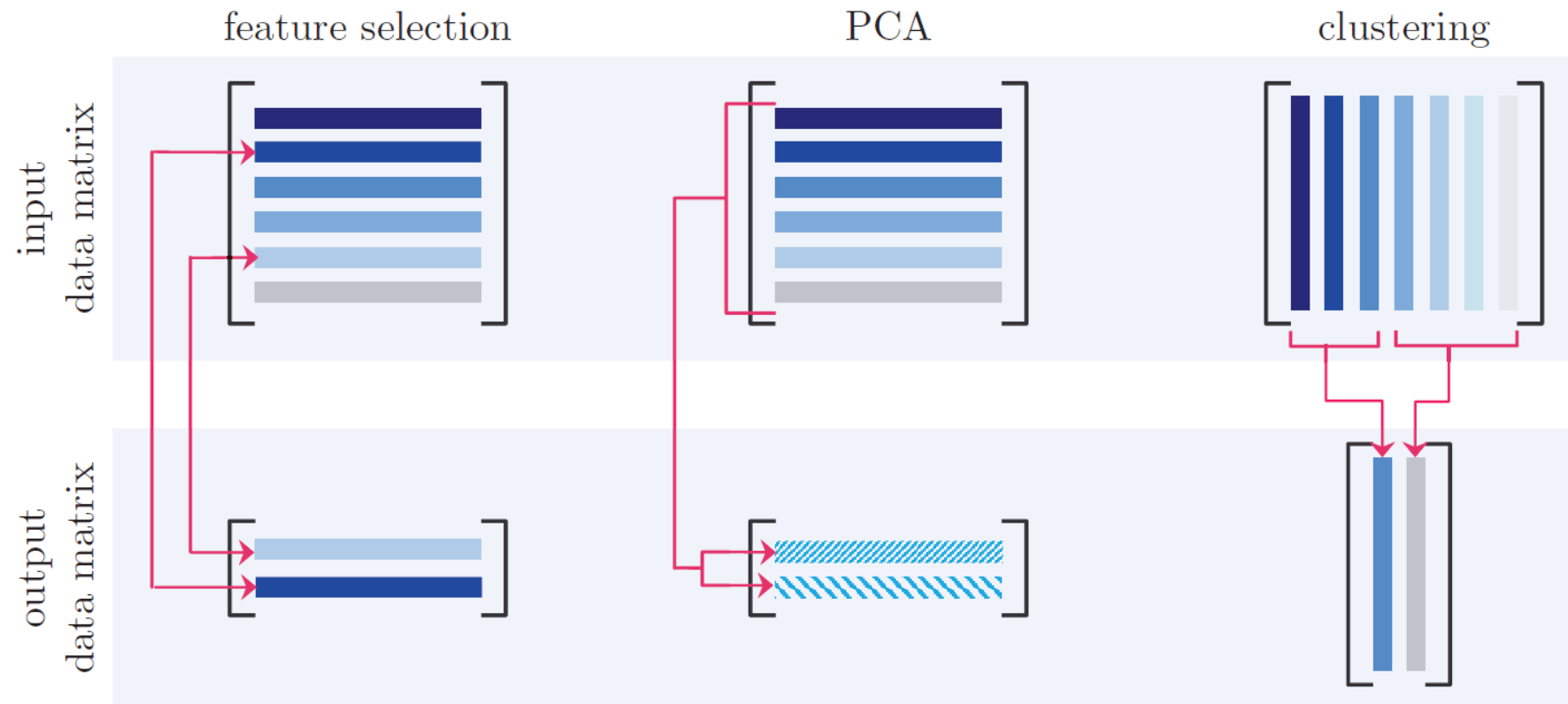
- Data dimension reduction

# Data dimension reduction

---

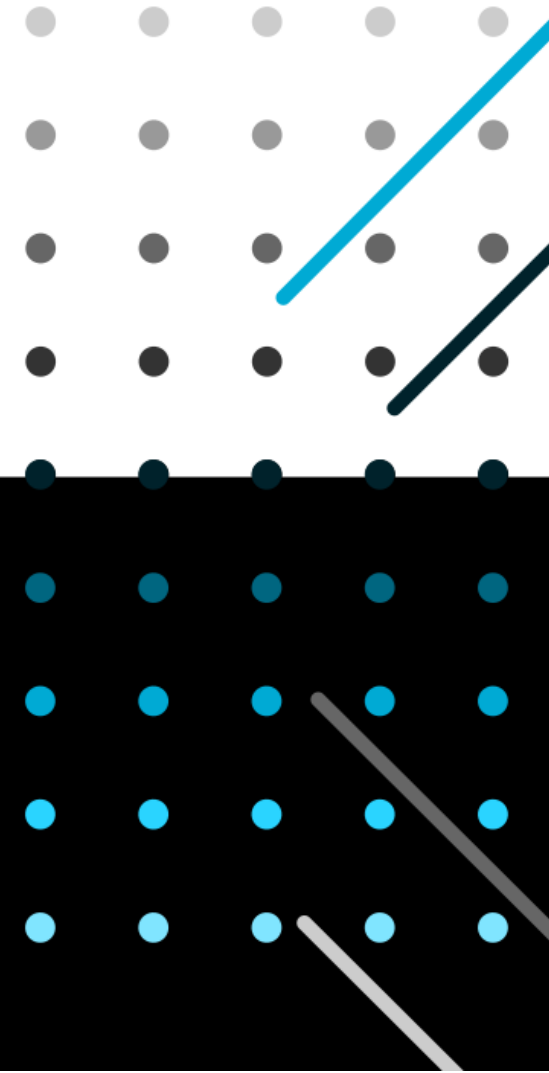


# Data dimension reduction



# PCA: Principal Components Analysis

---

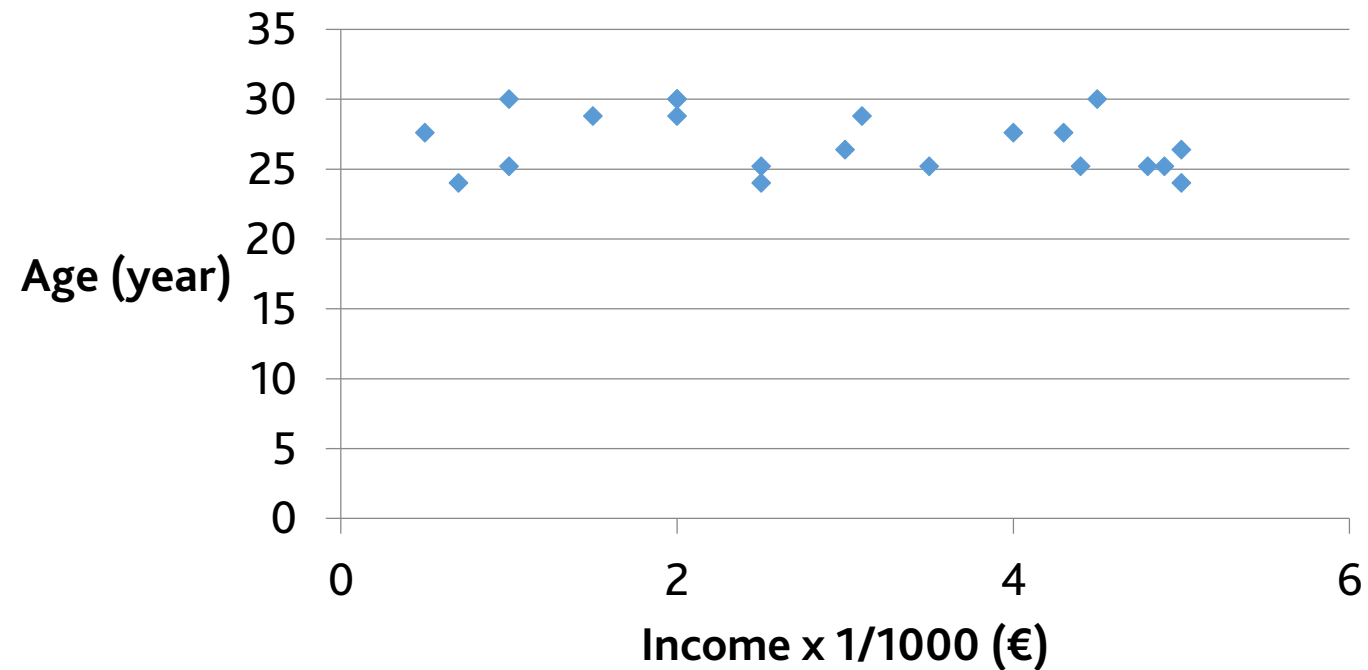


- What is PCA?
  - Linear transformation, invented in 1901 by Pearson and in 1933 by Hotteling.
  - Known also as « Principal Factor Analysis », referring to its first application in 1954 by Goodall.
  - Summarizes correlated data of  $l$  attributes by  $K$  uncorrelated axes (Principal Components).
  - The first component displays the maximum variance, the second displays the second maximum variance and so on.
- Why PCA?
  - Efficient reduction of the data.
  - Better visualization of the data.
  - Better classification.
  - Etc.

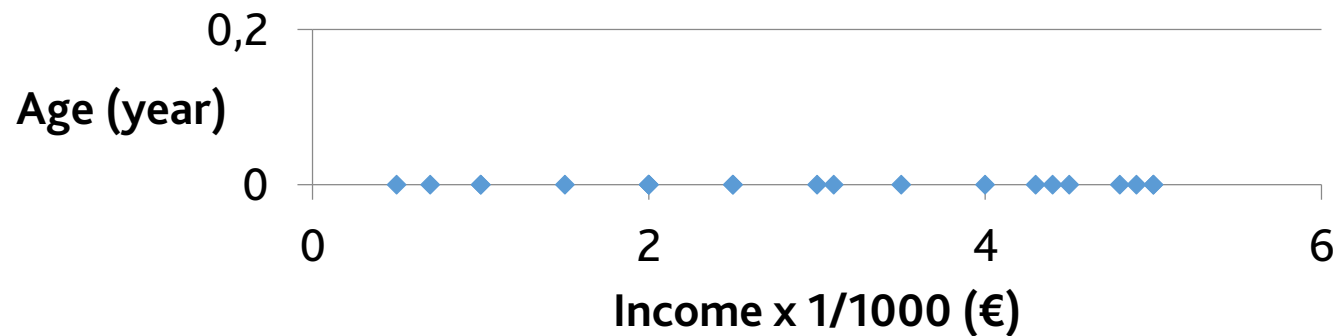
- Examples:
  - How to efficiently present a data  $X$  of size  $l \times N$  by  $K$  principal components without information loss?
    - $X: \mathbf{x}_1, \dots, \mathbf{x}_N \rightarrow X': \mathbf{x}'_1, \dots, \mathbf{x}'_N: \mathbb{R}^l \rightarrow \mathbb{R}^K$
  - How to achieve the same or better accuracy with less dimensions?
    - Classification
    - Clustering
    - Any Machine learning / data mining task



- Overview:
  - Considering the following samples, each of which is presented by income and age.
  - Are both attributes important to understand the data?

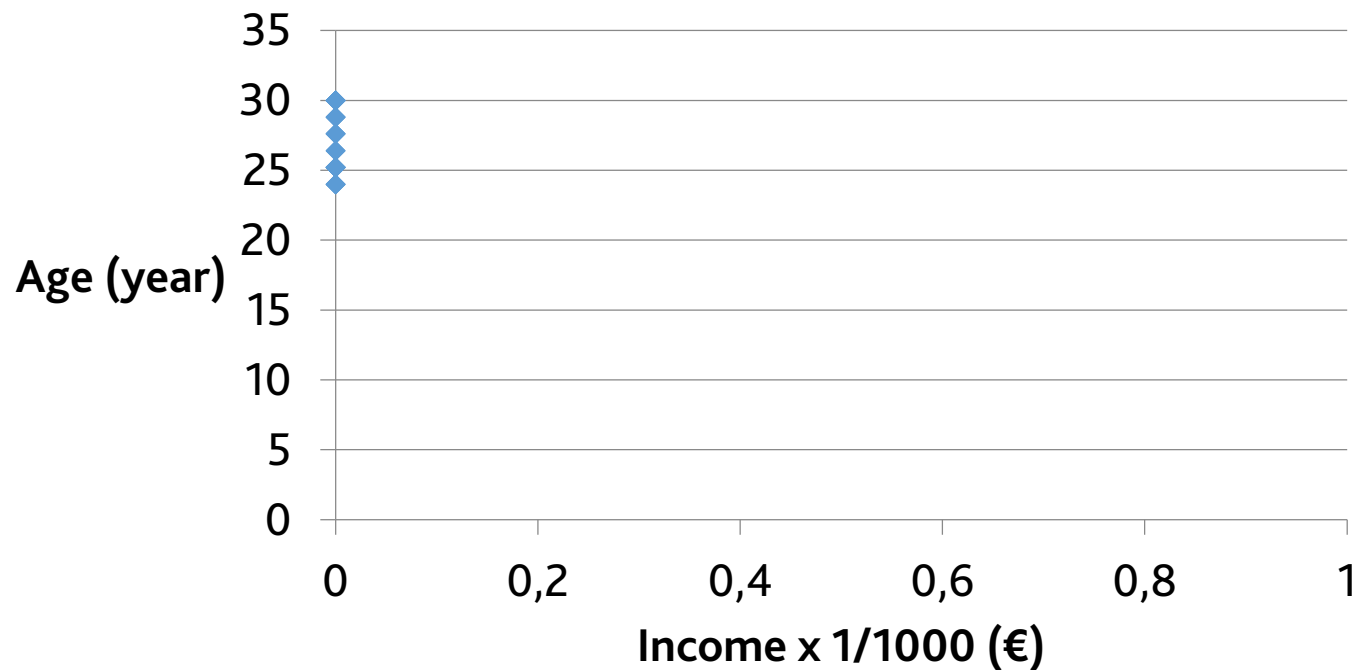


- Overview:
  - What if we simply ignore age, does this change anything in understanding the data?
  - Age does not give such valuable information which can help to understand and analyze the data.

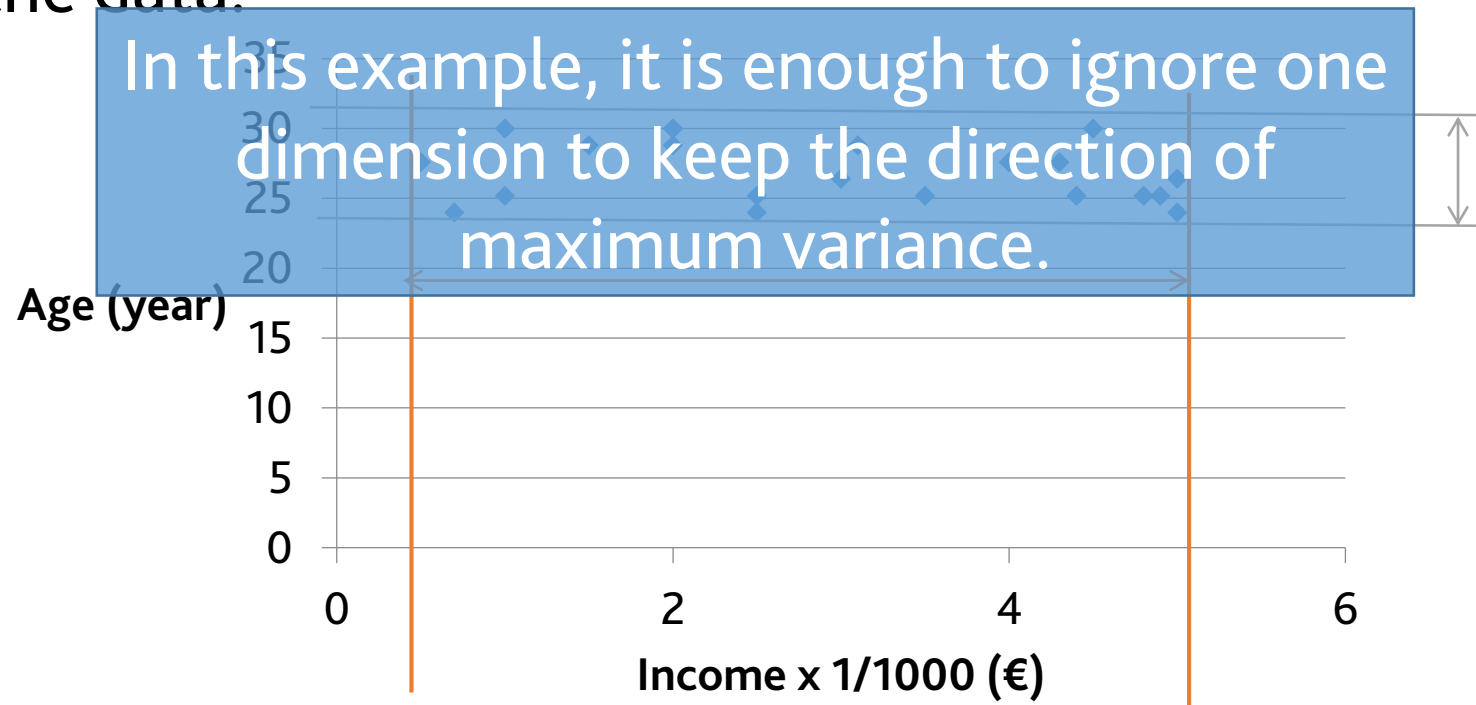


# PCA: Principal Components Analysis

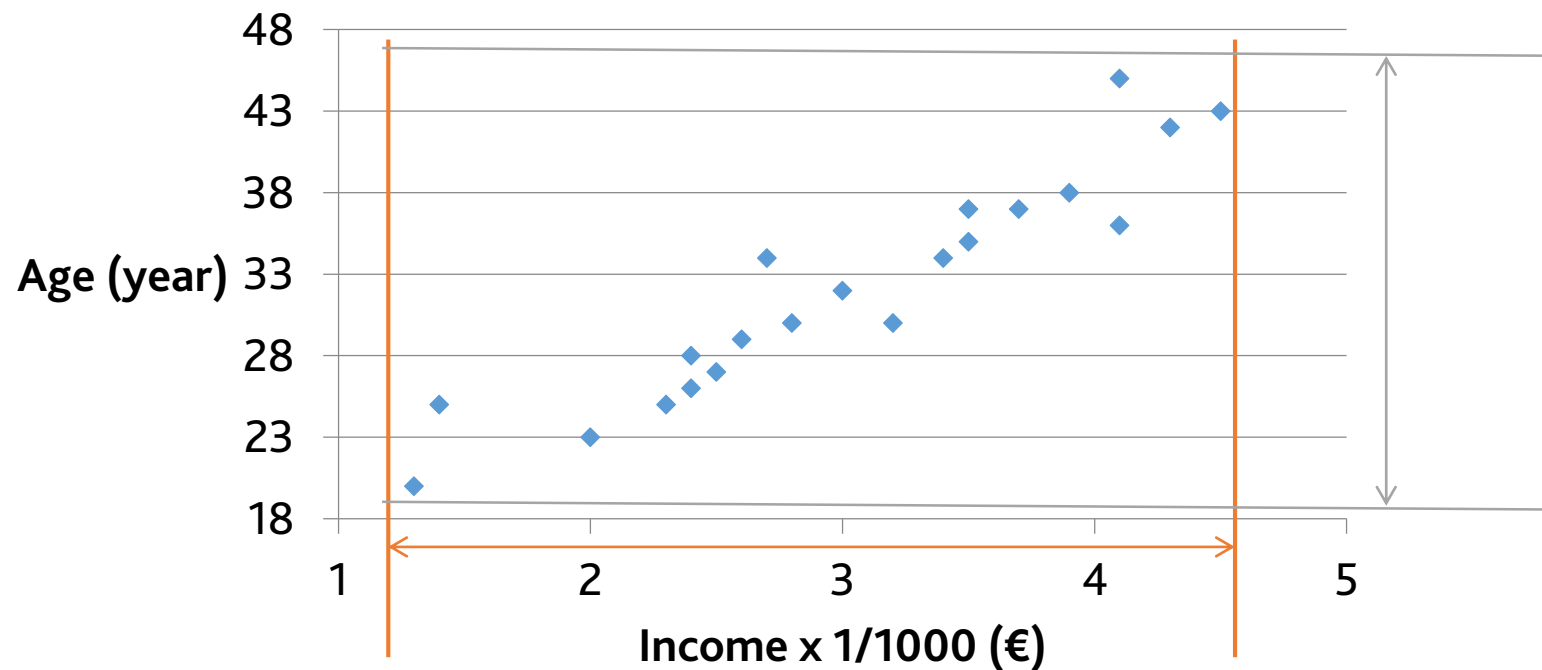
- Overview:
  - What if we simply ignore income, does this change anything in understanding the data?
  - Income was important to keep a high variance among samples.
  - But, what makes us consider Income much important feature than age by only visual assessments?



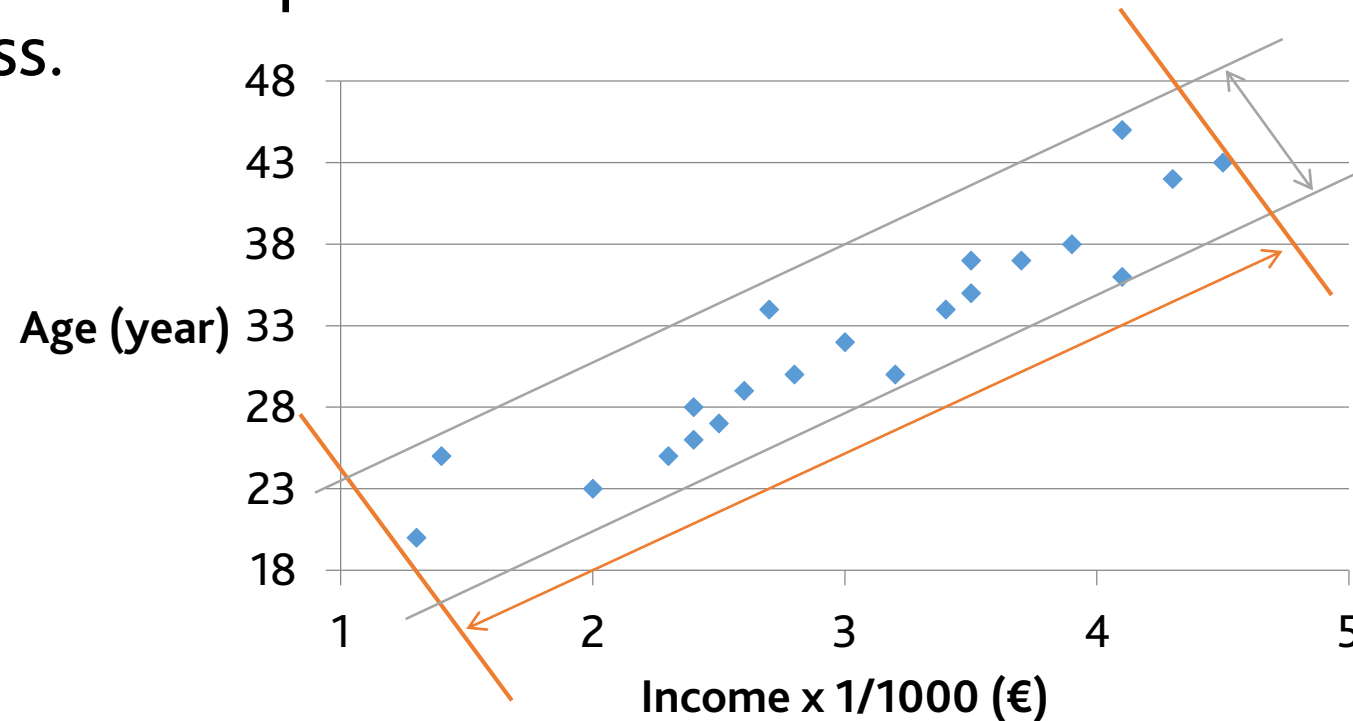
- Overview:
  - what makes us consider Income much important feature than age by only visual assessment?
  - **The answer is:** how the data is spread out or in another word the variance of the data.



- Overview:
  - It is not possible to flat the data on one axis because the one-d variations along both dimensions are quite similar.
  - However, if we consider a rotated coordinate system ( $45^\circ$ ), we can make it different.

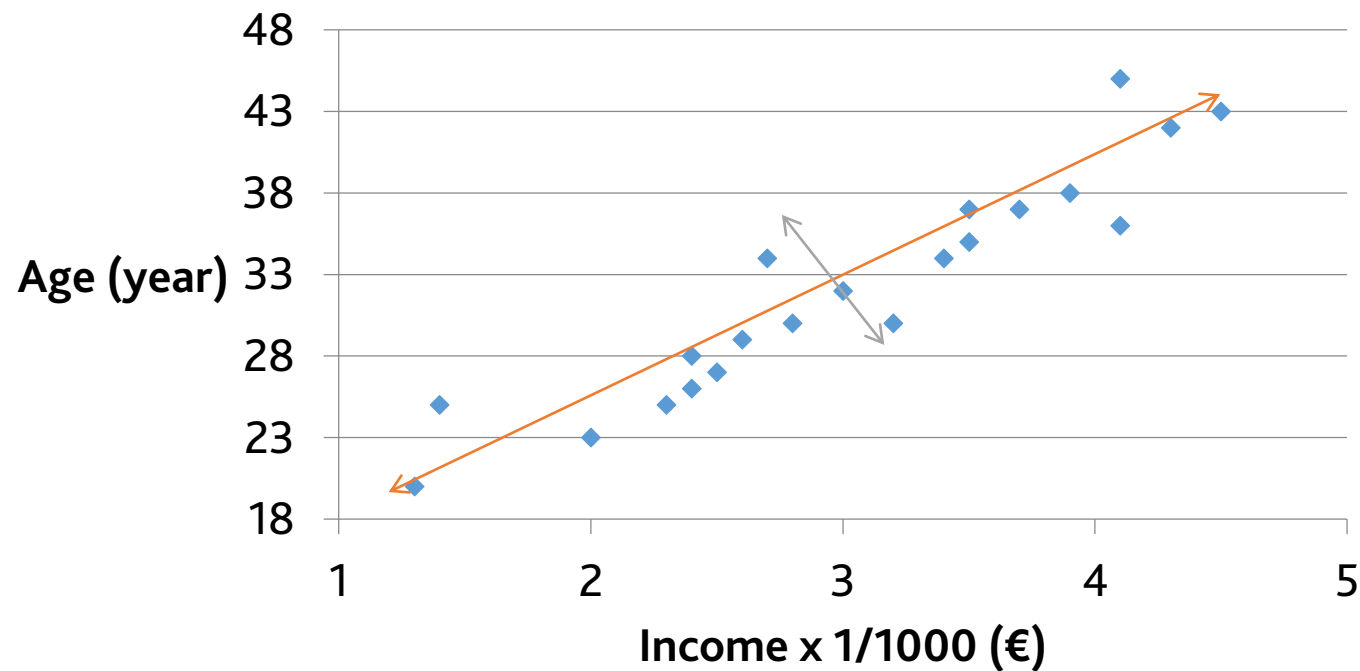


- Overview:
  - The new axes in the rotated coordinate system represent neither age nor income, but the combination of both.
  - Now, the data can be represented on one axis without too much information loss.

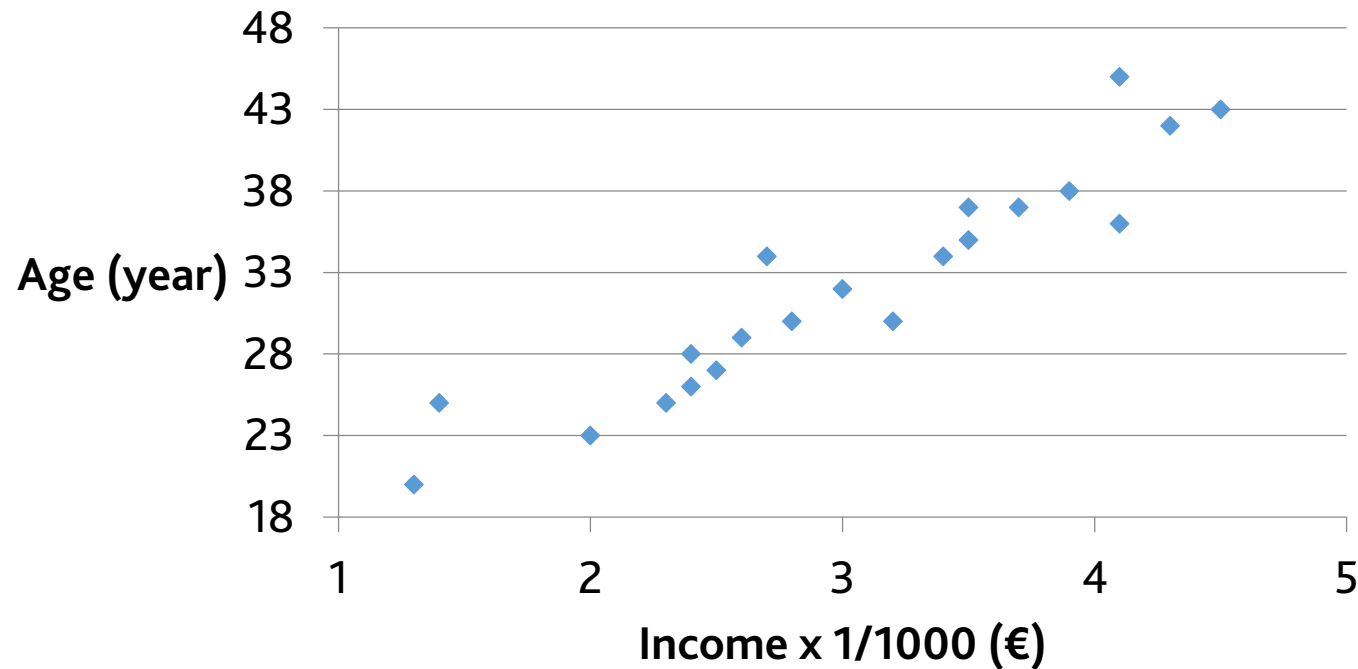


# PCA: Principal Components Analysis

- The new axes that span the direction of the highest variances are called Principal components.
- Note that the first PC (Principal Component) describes the highest variance, the second PC describes the second highest variance and the  $k^{th}$  PC describes the  $k^{th}$  highest variance.



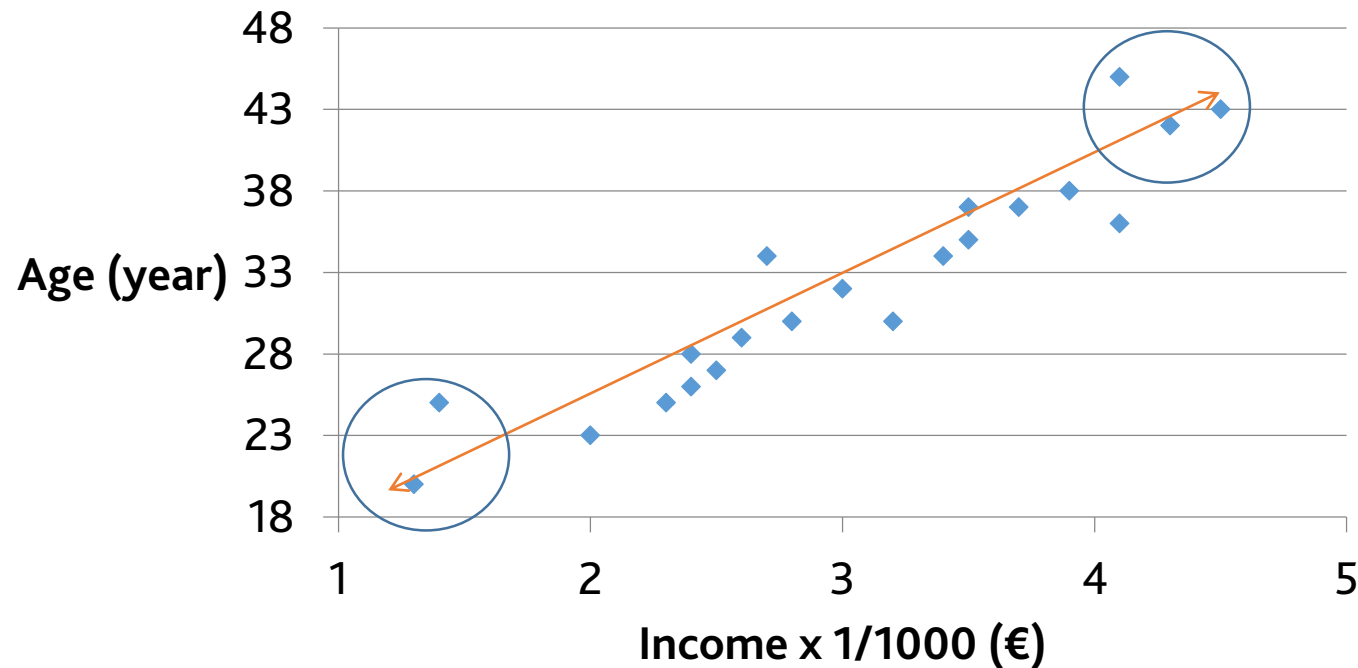
- In this example, which data points influence more the PCs in terms of length and direction?
- Considering only PC1:





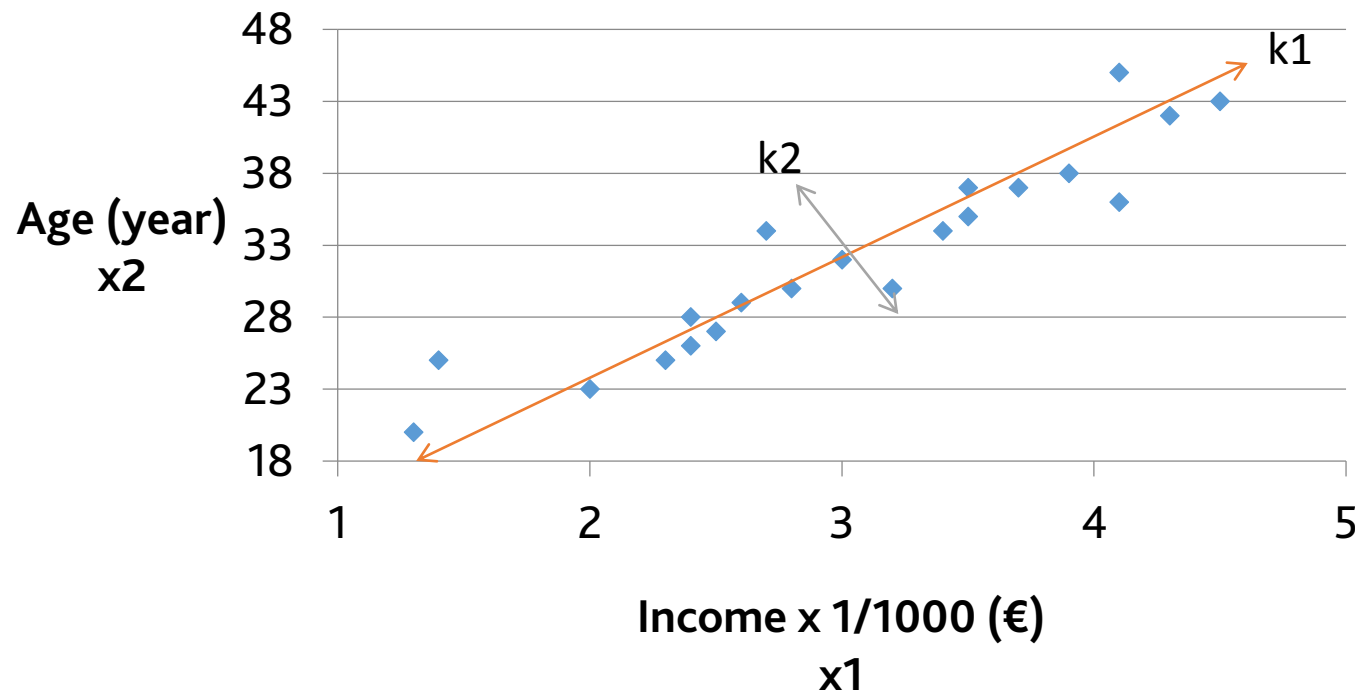
# PCA: Principal Components Analysis

- In this example, which datapoints influence more the PCs in terms of length and direction?
- Considering only PC1:
- The ones at endpoints influence more the PC, since they maximize the variance more than other data points.



# PCA: Principal Components Analysis

- Given  $N$  training samples:  $D: \{(x_i, y_i)\}_{i=1}^N$ , where  $x_i$  and  $y_i$  denotes the feature vector and the label of the  $i$ th instance, respectively.
- **Goal:** Project  $X$  of dimension  $l \times N$  onto  $X'$  having dimensionality  $K \times N (K < l)$  while maximizing the variance of  $X'$ .



# PCA: Principal Components Analysis

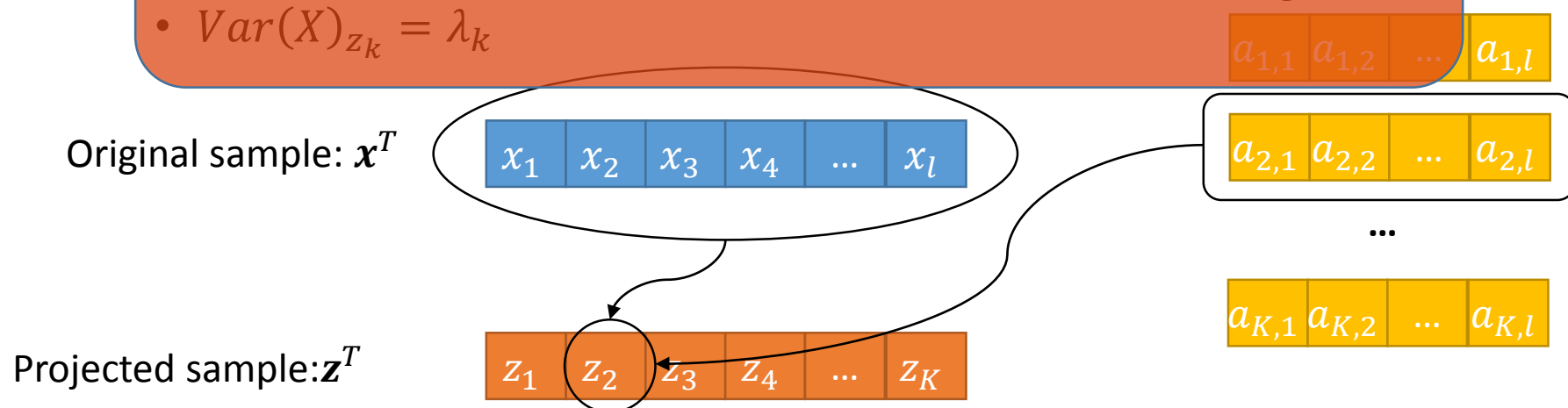
- Goal: find the  $K$  principal vectors such that:

$$z_k = \mathbf{a}_k^T \mathbf{x} = \sum_{j=1}^l a_{j,k} x_j$$

- Where

- $z_k$  is the  $k$ th principal component.
  - $|\mathbf{a}_k| = l$ , for  $k = 1, 2, \dots, K$
  - $\mathbf{a}_1$  ensures the maximum variance.
  - $\mathbf{a}_k$  ensures the  $k$ th maximum variance.
  - $Var(X)_{z_k} = \lambda_k$
- The order is controlled by the eigenvalues

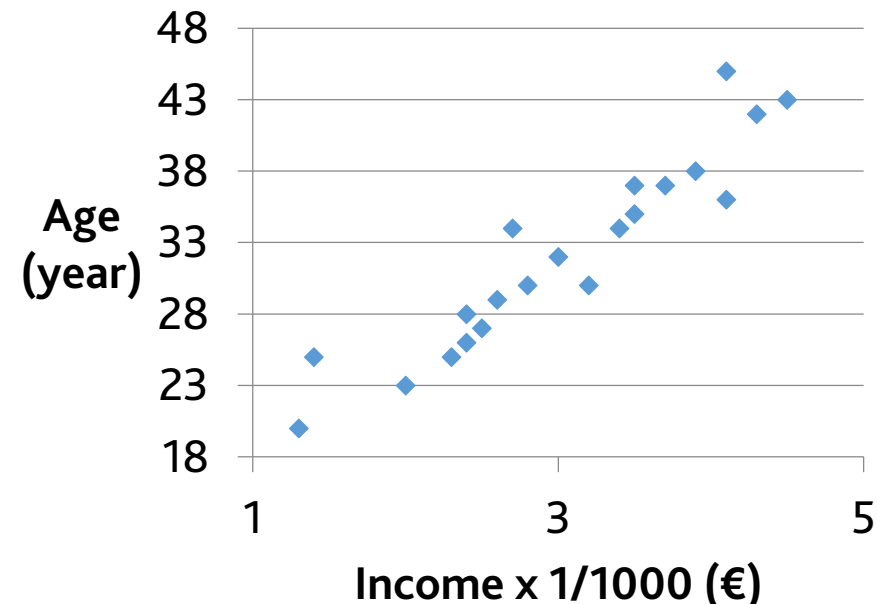
$K$  Eigenvectors

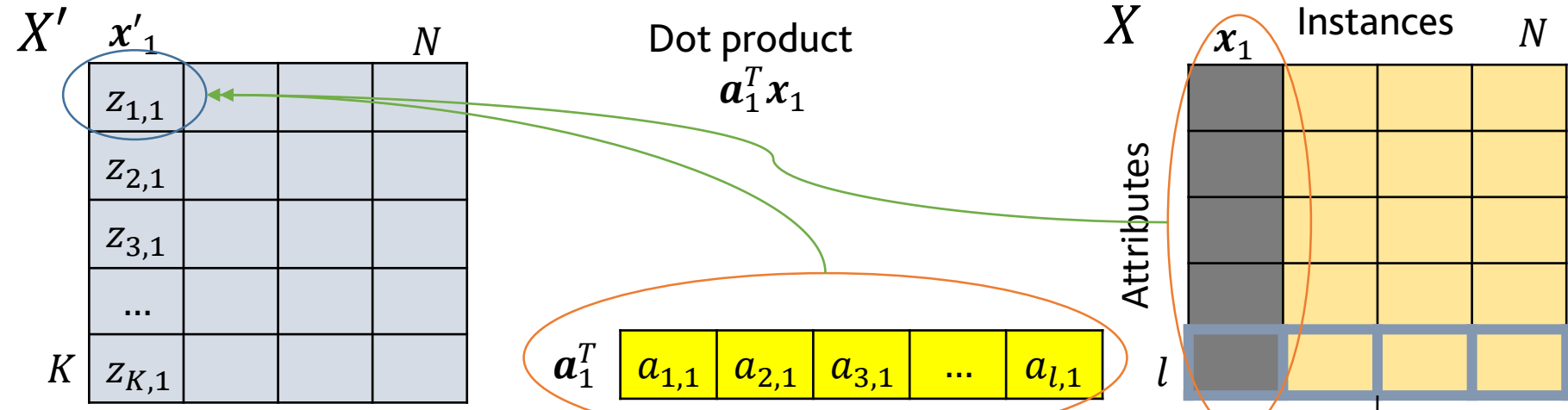


- Goal: find  $\mathbf{a}_{k=1}$  that maximizes  $Var(\mathbf{z}_{k=1})$ 
  - Until  $k \rightarrow K$
- $\mathbf{a}_k$  may contain very big values
  - $\mathbf{a}_k^T \mathbf{a}_k = 1$  (unit vector)

- $Var(income) = \frac{1}{N} \sum_{i=1}^N (x_{i,income} - \bar{x}_{income})^2$
- Given that  $z_k = \mathbf{a}_k^T \mathbf{x}$ :
  - $Var(\mathbf{z}_k) = \frac{1}{N} \sum_{i=1}^N (\mathbf{a}_k^T \mathbf{x}_i - \mathbf{a}_k^T \bar{X})^2$ 
    - $\bar{X}$  is the mean over all samples
  - $Var(\mathbf{z}_k) = \mathbf{a}_k^T \mathbf{S} \mathbf{a}_k$ 
    - $\mathbf{S}$  is the covariance matrix
  - $\mathbf{S} = Cov(X)$
  - $= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{X})(\mathbf{x}_i - \bar{X})^T$

See next « bonus » slide!

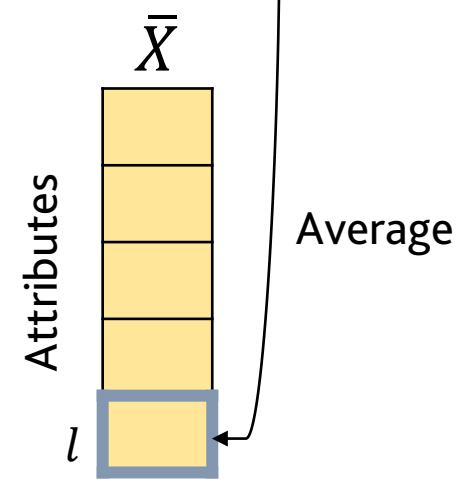




We ignore the second index in other slides for the sake of simplicity.

Why  $\overline{\mathbf{a}_k^T X} = \mathbf{a}_k^T \bar{X}$  ?

$$\begin{aligned}
 \overline{\mathbf{a}_k^T X} &= \frac{1}{N} \sum_{i=1}^N \mathbf{a}_k^T \mathbf{x}_i \\
 &= \frac{\mathbf{a}_k^T}{N} \sum_{i=1}^N \mathbf{x}_i \\
 &= \mathbf{a}_k^T \bar{X}
 \end{aligned}$$



- For the first PC, the task becomes now:
  - Maximizing  $\mathbf{a}_1^T \mathbf{S} \mathbf{a}_1$  with the constraint  $\mathbf{a}_k^T \mathbf{a}_k = 1$ 
    - Lagrange multiplier

- The maximization becomes unconstrained
$$\mathbf{a}_1^T \mathbf{S} \mathbf{a}_1 + \lambda_1 (1 - \mathbf{a}_1^T \mathbf{a}_1)$$

- By differentiating w.r.t  $\mathbf{a}_1$ 
$$\frac{\partial}{\partial \mathbf{a}_1} (\mathbf{a}_1^T \mathbf{S} \mathbf{a}_1 - \lambda_1 (\mathbf{a}_1^T \mathbf{a}_1 - 1)) = 0$$
$$\Rightarrow \mathbf{S} \mathbf{a}_1 - \lambda_1 \mathbf{a}_1 = 0$$
$$\Rightarrow \mathbf{S} \mathbf{a}_1 = \lambda_1 \mathbf{a}_1$$
$$\Rightarrow \mathbf{a}_1^T \mathbf{S} \mathbf{a}_1 = \lambda_1$$

Remember:  $\mathbf{a}_1^T \mathbf{a}_1 = 1$

The variance will be a maximum when we set  $\mathbf{a}_1$  equals to the eigenvector having the largest eigenvalue  $\lambda_1$

- $\mathbf{a}_1$  is called the first principal component
- The second principal component is the one that maximises  $\mathbf{a}_2^T \mathbf{S} \mathbf{a}_2$  subject to:  $\mathbf{z}_2$  and  $\mathbf{z}_1$  are uncorrelated.

- $Cov(\mathbf{z}_1, \mathbf{z}_2) = 0$

$$\Rightarrow \mathbf{a}_1^T \mathbf{S} \mathbf{a}_2 = 0$$

$$\Rightarrow \mathbf{a}_2^T \mathbf{S} \mathbf{a}_1 = 0$$

$$\Rightarrow \mathbf{a}_2^T \lambda_1 \mathbf{a}_1 = 0$$

- Two constraints:

- $\mathbf{a}_2^T \mathbf{a}_2 = 1$

- $\mathbf{a}_2^T \lambda_1 \mathbf{a}_1 = 0$

- Two Lagrange multipliers:  $\lambda_2$  and  $\phi$





- The unconstrained maximization of  $\mathbf{a}_2^T \mathbf{S} \mathbf{a}_2$  becomes
$$\mathbf{a}_2^T \mathbf{S} \mathbf{a}_2 + \lambda_2 (1 - \mathbf{a}_2^T \mathbf{a}_2) + \phi (0 - \mathbf{a}_2^T \mathbf{a}_1)$$

- By differentiating w.r.t  $\mathbf{a}_2$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{a}_2} (\mathbf{a}_2^T \mathbf{S} \mathbf{a}_2 - \lambda_2 (\mathbf{a}_2^T \mathbf{a}_2 - 1) - \phi \mathbf{a}_2^T \mathbf{a}_1) &= 0 \\ \Rightarrow \mathbf{S} \mathbf{a}_2 - \lambda_2 \mathbf{a}_2^T - \phi \mathbf{a}_1 &= 0 \\ \Rightarrow \mathbf{a}_1^T \mathbf{S} \mathbf{a}_2 - \lambda_2 \mathbf{a}_1^T \mathbf{a}_2^T - \phi \mathbf{a}_1^T \mathbf{a}_1 &= 0 \\ \Rightarrow 0 - 0 - \phi &= 0 \\ \Rightarrow \phi &= 0 \\ \Rightarrow \mathbf{S} \mathbf{a}_2 - \lambda_2 \mathbf{a}_2^T &= 0 \\ \Rightarrow \mathbf{a}_2^T \mathbf{S} \mathbf{a}_2 &= \lambda_2 \end{aligned}$$

- $\mathbf{a}_2$  is the eigenvector associated with the second largest eigenvalue  $\lambda_2$  yielding the second PC
- This process can be repeated up to  $K = l$  eigenvectors.

- Let  $A$  be an orthogonal  $K \times l$  matrix consisting of  $K$  eigenvectors:
  - $\mathbf{z} = A^T \mathbf{x}$
  - Then:  $Cov(Z) = \Lambda = A^T \mathbf{S} A$
  - For the data matrix  $X$ :
    - $Z = A^T X$

- How to compute the eigenvectors and eigenvalues of  $S$ ?
  - If  $S$  is a square matrix, a non zero vector  $\mathbf{a}_k$  is an eigenvector of  $S$  such that:  $S\mathbf{a}_k = \lambda_k \mathbf{a}_k$ , where  $\lambda_k$  is the corresponding eigenvalue. For example:
    - $\begin{pmatrix} 2 & 6 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 6 \\ 2 \end{pmatrix} = \begin{pmatrix} 24 \\ 8 \end{pmatrix} = 4 \begin{pmatrix} 6 \\ 2 \end{pmatrix}$
  - Considering all eigenvectors:

$$SA = \lambda IA$$

Let  $\Psi = \lambda I$

$$S = A\Psi A^T$$

Remember  $A$  is orthonormal  $\Rightarrow \mathbf{a}_k^T \mathbf{a}_k = 1 \Rightarrow A^{-1} = A^T$

## Eigendecomposition

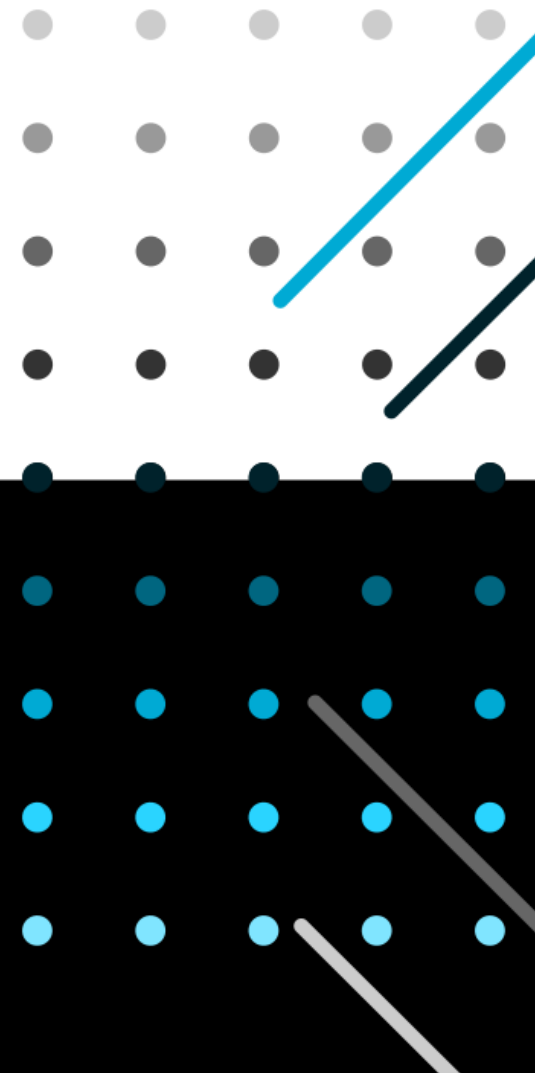
See <https://www.scss.tcd.ie/Rozenn.Dahyot/CS1BA1/SolutionEigen.pdf>



- Data standardization:  $X = \text{standardization}(X)$
- Covariance matrix calculation:  $S = \text{Cov}(X)$
- Finding eigenvectors and eigenvalues
- Eigenvectors ordering (descending) w.r.t eigenvalues
- Project  $X$  w.r.t the eigenvectors.

# SVD: Singular Value Decomposition

---



# SVD: Singular Value Decomposition

- Given that any  $a \times b$  matrix  $X$  can be uniquely expressed as:

$$X_{a \times b} = U_{a \times a} \Sigma_{a \times b} V_{b \times b}^T$$

- Where:

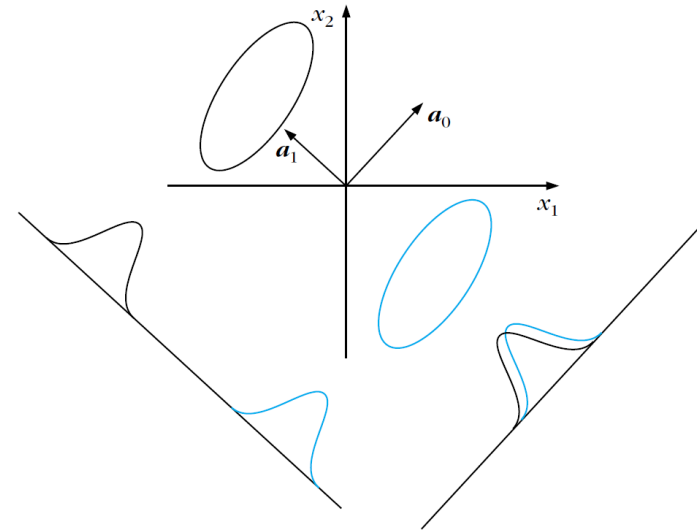
- $U$  is a unity matrix
- $\Sigma$  is the diagonal matrix of singular values.

- Assuming that  $X$  is centred, i.e. zero mean:

$$\begin{aligned} \text{Cov}(X) &= \frac{1}{N-1} X^T X \\ &= \frac{1}{N-1} V \Sigma U^T U \Sigma V^T \\ &= V \frac{\Sigma^2}{N-1} V^T \end{aligned}$$

- Given that
  - $Cov(X) = \mathbf{S} = A\Psi A^T$
  - $Cov(X) = V \frac{\Sigma^2}{N-1} V^T$
- This means that  $V$  are principal vectors (eigenvectors) and the diagonal of  $\frac{\Sigma^2}{N-1}$  are eigenvalues
- **Very important: This is correct only when  $X$  is centred.**

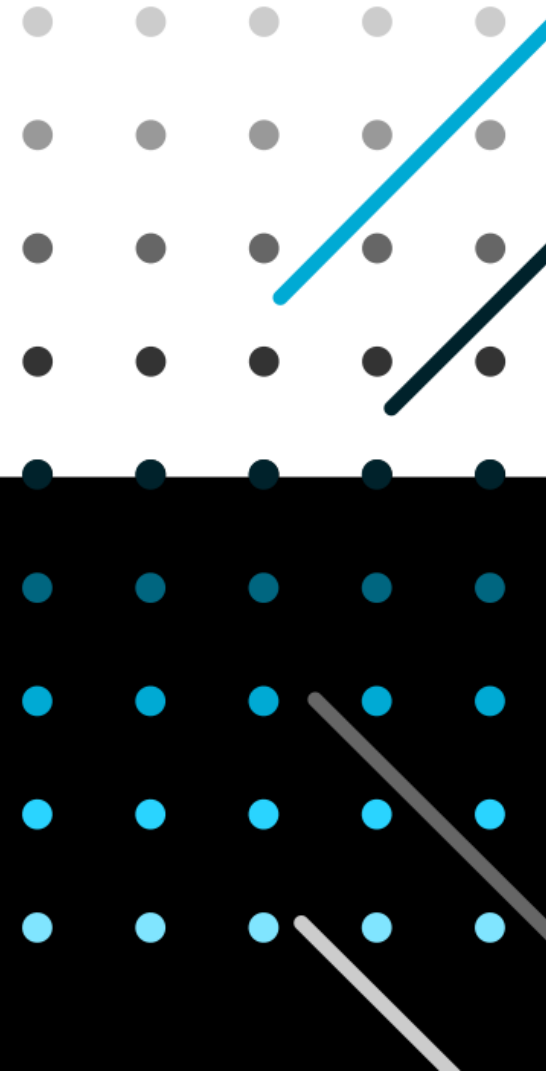
- Obviously, when we apply PCA on the training data, we need to use the same eigenvectors to transform the unseen data.
- Is PCA always helpful?
  - No, the eigenvector with the largest eigenvalue might make the two classes coincide. The classes might be better separated by another eigenvector.





# Summary

---



- Data dimension reduction
  - PCA
  - SVD

Thank you!



## Zeyd Boukhers

E-mail: [Boukhers@uni-koblenz.de](mailto:Boukhers@uni-koblenz.de)

Phone: +49 (0) 261 287-2765

Web: [Zeyd.Boukhers.com](http://Zeyd.Boukhers.com)

University of Koblenz-Landau

Universitätsstr. 1

56070 Koblenz

