

Lecture 11: PCA

Andreas Wichert
Department of Computer Science and Engineering
Técnico Lisboa

1

The Karhunen-Loève Transform

- The Karhunen-Loève transform is a linear transform that maps possibly correlated variables into a set of values of linearly uncorrelated variables
- This transformation is defined in such a way that the first principal component has the largest possible variance

2

The Covariance

- The sample size denoted by n , is the number of data items in a sample of a population.
- The goal is to make inferences about a population from a sample.
- Sample covariance indicates the relationship between two variables of a sample

$$cov(X, Y) = \frac{\sum_{k=1}^n (x_k - \bar{x}) \cdot (y_k - \bar{y})}{n - 1}$$

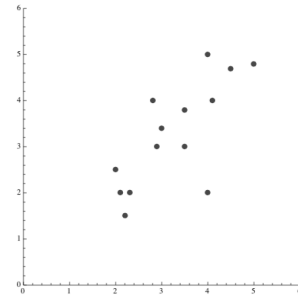
- The sample covariance has $n - 1$ in the denominator rather than n due to Bessel's correction.
- For the whole population, the covariance is

$$cov(X, Y) = \frac{\sum_{k=1}^n (x_k - \bar{x}) \cdot (y_k - \bar{y})}{n}.$$

3

- The sample covariance relies on the difference between each observation and the sample mean.
- In computer science, the sample covariance is usually used and
- In statistics (Bishop) population covariance
- In a linear relationship, either the high values of one variable are paired with the high values of another variable or the high values of one variable are paired with the low values of another variable

4



- For example, for a list of two variables, (X, Y) ,

$$\Sigma = \{(2.1, 2), (2.3, 2), (2.9, 3), (4.1, 4), (5, 4.8), (2, 2.5), (2.2, 1.5), \\ (4, 5), (4, 2), (2.8, 4), (3, 3.4), (3.5, 3.8), (4.5, 4.7), (3.5, 3)\}$$

- represents the data set Σ . The sample covariance of the data set is 0.82456 . Ordering the list by X , we notice that the **ascending** X values are matched by **ascending** Y values

5

- The covariance matrix measures the tendency of two features, x_i and x_j , to vary in the same direction. The covariance between features x_i and x_j is estimated for n vectors as

$$c_{ij} = \frac{\sum_{k=1}^n (x_{k,i} - \bar{x}_i) \cdot (y_{k,j} - \bar{y}_j)}{n - 1} \quad c_{ij} = \frac{\sum_{k=1}^n (x_i^{(k)} - m_i)(x_j^{(k)} - m_j)}{n - 1}$$

- with m_i and m_j being the arithmetic mean of the two variables of the sample. Covariances are symmetric; $c_{ij} = c_{ji}$

6

Correlation

- Covariance is related to correlation

$$r_{ij} = \frac{\sum_{k=1}^n (x_i^{(k)} - m_i)(x_j^{(k)} - m_j)}{(n-1)s_i s_j} = \frac{c_{ij}}{s_i s_j} \in [-1, 1]$$

7

- The resulting covariance matrix C is symmetric and positive-definite,

$$C = \begin{pmatrix} c_{11} & c_{12} & \cdots & c_{1m} \\ c_{21} & c_{22} & \cdots & c_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ c_{m1} & c_{m2} & \cdots & c_{mm} \end{pmatrix}$$

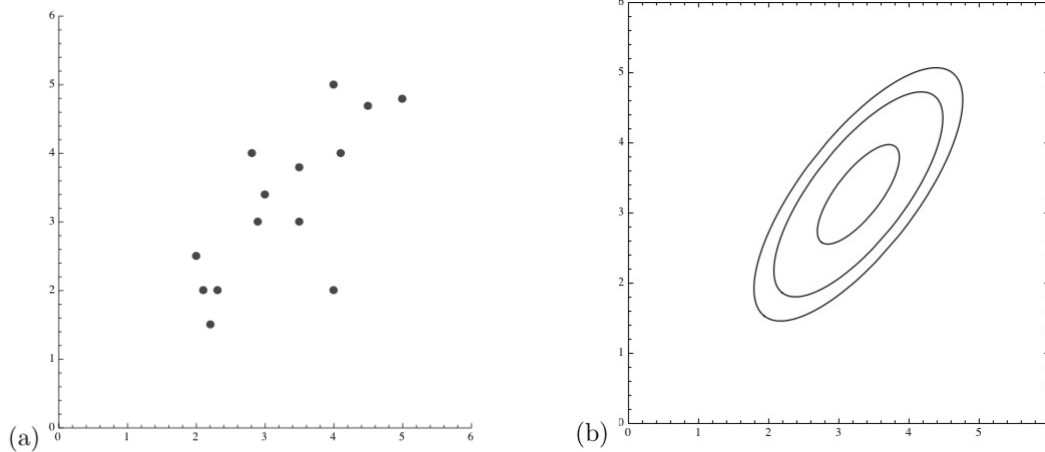
8

$$\Sigma = \{(2.1, 2), (2.3, 2), (2.9, 3), (4.1, 4), (5, 4.8), (2, 2.5), (2.2, 1.5), (4, 5), (4, 2), (2.8, 4), (3, 3.4), (3.5, 3.8), (4.5, 4.7), (3.5, 3)\}$$

$$c_{ij} = \frac{\sum_{k=1}^n (x_{k,i} - \bar{x}_i) \cdot (y_{k,j} - \bar{y}_j)}{n - 1}$$

$$C = \begin{pmatrix} 0.912582 & 0.82456 \\ 0.82456 & 1.34247 \end{pmatrix}$$

9



- (a) The data points of the set Σ (b) The two dimensional distribution of Σ can be described by three ellipse that divide the data points in four equal groups.

10

The Karhunen-Loève Transform

A real matrix M is positive definite if $\mathbf{z}^\top \cdot M \cdot \mathbf{z}$ is positive for any non-zero column vector \mathbf{z} of real numbers. A symmetric and positive-definite matrix can be diagonalized. It follows that

$$U^{-1} \cdot C \cdot U = \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$$

U is an orthonormal matrix of the dimension $m \times m$,

$$U^\top \cdot U = I$$

$$U^\top \cdot C \cdot U = \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$$

and

$$U \cdot \Lambda = C \cdot U.$$

11

$$U \cdot \Lambda = C \cdot U.$$

There are m eigenvalues and eigenvectors with

$$(\lambda_i \cdot I - C) \cdot \mathbf{u}_i = 0$$

and

$$C \cdot \mathbf{u}_i = \lambda_i \cdot \mathbf{u}_i$$

An eigenvector can have two directions, it is either \mathbf{u}_i or $-\mathbf{u}_i$.

$$C \cdot (-\mathbf{u}_i) = \lambda_i \cdot (-\mathbf{u}_i)$$

12

The eigenvectors are always orthogonal, and their length is arbitrary. The normalized eigenvectors define the orthonormal matrix U of dimension $m \times m$. Each normalized eigenvector is a column of U with

$$U^T \cdot U = I.$$

The matrix U defines the Karhunen-Loève transform. The Karhunen-Loève transform rotates the coordinate system in such a way that the new covariance matrix will be diagonal

$$\mathbf{y} = U^T \cdot \mathbf{x}$$

13

- The squares of the eigenvalues represent the variances along the eigenvectors. The eigenvalues corresponding to the covariance matrix of the data set Σ are

$$\lambda_1 = 1.97964, \quad \lambda_2 = 0.275412$$

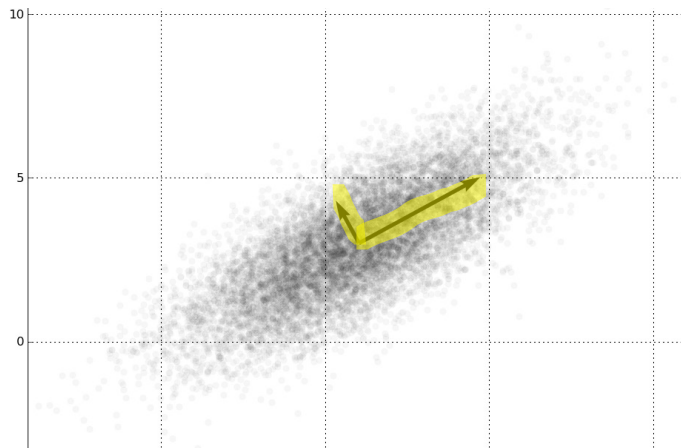
and the corresponding normalized eigenvectors are

$$\mathbf{u}_1 = \begin{pmatrix} 0.611454 \\ 0.79128 \end{pmatrix}, \quad \mathbf{u}_2 = \begin{pmatrix} -0.79128 \\ 0.611454 \end{pmatrix}.$$

We define the matrix U with

$$U = \begin{pmatrix} 0.611454 & -0.79128 \\ 0.79128 & 0.611454 \end{pmatrix}.$$

14



15

The define the matrix U with

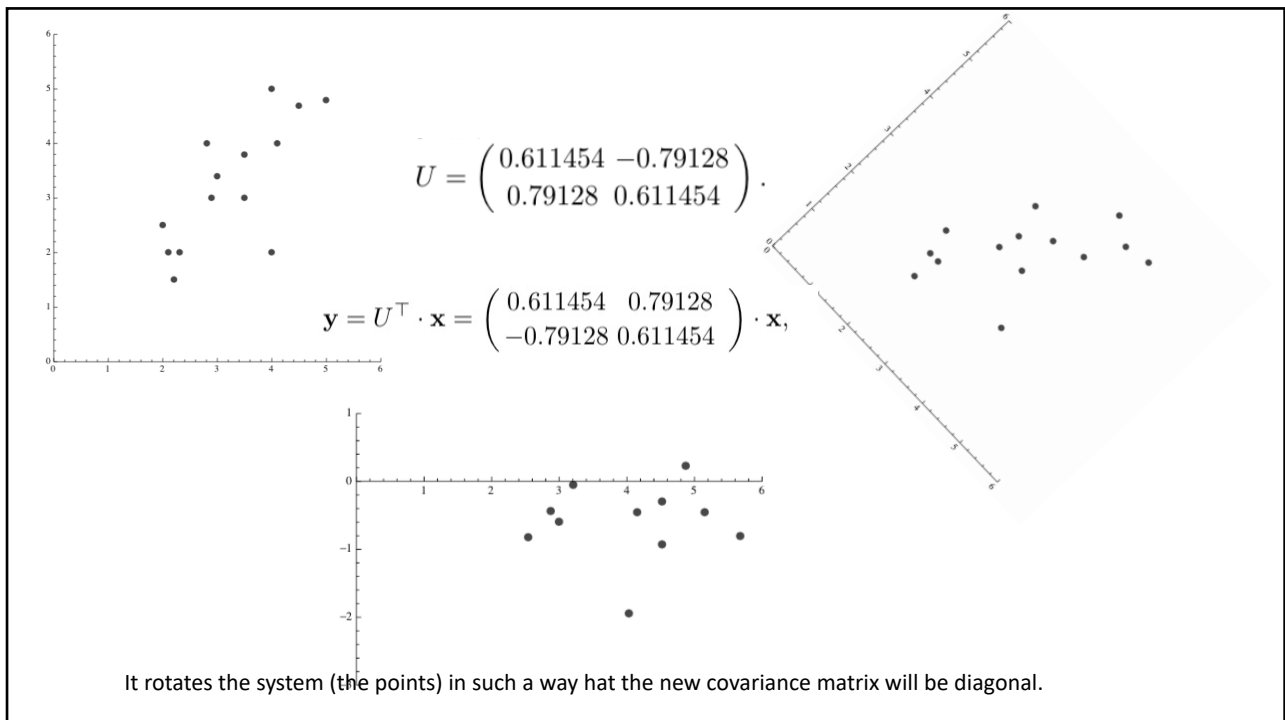
$$U = \begin{pmatrix} 0.611454 & -0.79128 \\ 0.79128 & 0.611454 \end{pmatrix}.$$

The Karhunen-Loève transform for the data set Σ is given by

$$\mathbf{y} = U^{\top} \cdot \mathbf{x} = \begin{pmatrix} 0.611454 & 0.79128 \\ -0.79128 & 0.611454 \end{pmatrix} \cdot \mathbf{x},$$

it rotates the coordinate system in such a way that the new covariance matrix will be diagonal

16



17

Principal component analysis

- Principal component analysis (PCA) is a technique that is useful for the compression of data.
- The purpose is to reduce the dimensionality of a data set by finding a new set of variables, smaller than the original set of variables, that nonetheless retains most of the sample's information.
- The first principal component corresponds to the normalized eigenvector with the highest variance.
- In principal component analysis (PCA), the significant eigenvectors define the principal components.

18

- Accordingly to the **Kaiser criterion**, the eigenvectors whose eigenvalues are below 1 are discarded
- Each of the s non-discarded eigenvectors is a column of the matrix W of dimension $s \times m$ with the linear mapping from $\mathbf{R}^m \rightarrow \mathbf{R}^s$,

$$\mathbf{z} = W^T \cdot \mathbf{x}$$

- The Principal component analysis for the data set \mathcal{Z} is given by

$$\mathbf{z} = W^T \cdot \mathbf{x} = (0.611454 \ 0.79128) \cdot \mathbf{x}$$

19

$$\lambda_1 = 1.97964, \quad \lambda_2 = 0.275412$$

and the corresponding normalized eigenvectors are

$$\mathbf{u}_1 = \begin{pmatrix} 0.611454 \\ 0.79128 \end{pmatrix}, \quad \mathbf{u}_2 = \begin{pmatrix} -0.79128 \\ 0.611454 \end{pmatrix}.$$

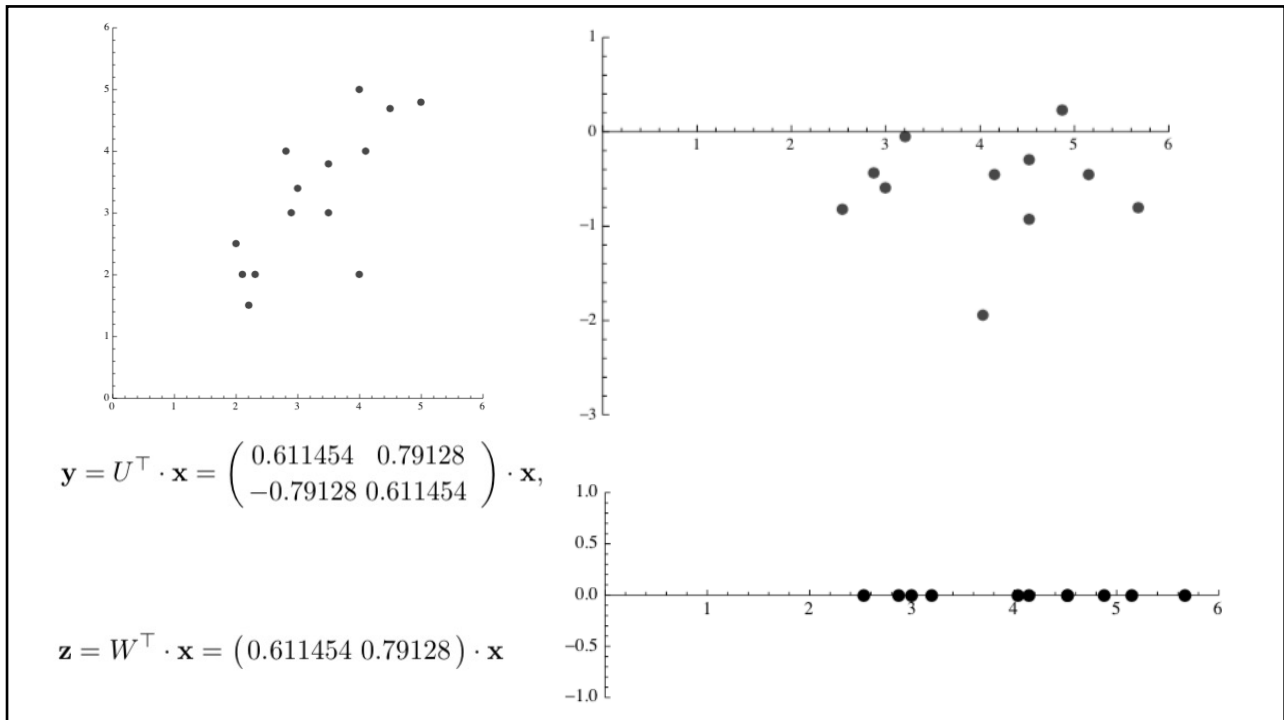
The define the matrix U with

$$U = \begin{pmatrix} 0.611454 & -0.79128 \\ 0.79128 & 0.611454 \end{pmatrix}.$$

$$\mathbf{y} = U^T \cdot \mathbf{x} = \begin{pmatrix} 0.611454 & 0.79128 \\ -0.79128 & 0.611454 \end{pmatrix} \cdot \mathbf{x},$$

$$\mathbf{z} = W^T \cdot \mathbf{x} = (0.611454 \ 0.79128) \cdot \mathbf{x}$$

20



21

- Suppose we have a covariance matrix

$$C = \begin{pmatrix} 3 & 1 \\ 1 & 21 \end{pmatrix}$$

- What is the corresponding matrix of the K-L transformation?
- First, we have to compute the eigenvalues.
- The system has to become linear depend-able (singular).
- The determinant has to become zero.

$$|\lambda \cdot I - C| = 0.$$

22

$$|\lambda \cdot I - C| = 0.$$

Solving the Equation

$$\lambda^2 - 24 \cdot \lambda + 62 = 0$$

we get the two eigenvalues

$$\lambda_1 = 2.94461, \quad \lambda_2 = 21.05538.$$

23

$$\lambda_1 = 2.94461, \quad \lambda_2 = 21.05538.$$

To compute the eigenvectors we have to solve two singular, dependent systems

$$|\lambda_1 \cdot I - C| = 0$$

and

$$|\lambda_2 \cdot I - C| = 0.$$

24

For $\lambda_1 = 2.94461$ we get

$$\left(\begin{pmatrix} 2.94461 & 0 \\ 0 & 2.94461 \end{pmatrix} - \begin{pmatrix} 3 & 1 \\ 1 & 21 \end{pmatrix} \right) \cdot \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = 0$$

and we have to find a nontrivial solution for

$$\begin{pmatrix} -0.05538 & -1 \\ -1 & -18.055 \end{pmatrix} \cdot \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = 0$$

Because the system is linear dependable, the left column is a multiple value of the right column, and there are infinitely many solution. We only have to determine the direction of the eigenvectors; if we simply suppose that $u_1 = 1$,

25

$$u_1 = 1,$$

$$\begin{pmatrix} -0.05538 & -1 \\ -1 & -18.055 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ u_2 \end{pmatrix} = 0$$

and

$$\begin{pmatrix} -0.05538 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ 18.055 \end{pmatrix} \cdot u_2$$

with

$$\mathbf{u}_1 = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} 1 \\ -0.05539 \end{pmatrix}.$$

26

For $\lambda_2 = 21.05538$ we get

$$\begin{pmatrix} 18.055 & -1 \\ -1 & 0.05538 \end{pmatrix} \cdot \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = 0$$

with

$$\mathbf{u}_2 = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 18.055 \end{pmatrix}.$$

The two normalized vectors $\mathbf{u}_1, \mathbf{u}_2$ define the columns of the matrix U

$$U = \begin{pmatrix} 0.998469 & 0.0553016 \\ -0.0553052 & 0.99847 \end{pmatrix}.$$

Because $\lambda_1 = 2.94461 < \lambda_2 = 21.05538$ the second eigenvector is more significant, however we can not apply the Kaiser criterion.

27

$$\Psi = \{(1, 1), (2, 2), (3, 3), (4, 4), (5, 5), (6, 6)\}$$

the covariance matrix is

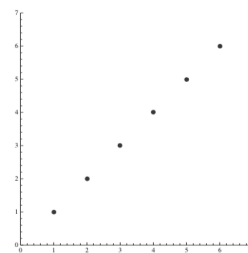
$$C = \begin{pmatrix} 3.5 & 3.5 \\ 3.5 & 3.5 \end{pmatrix}.$$

The two eigenvalues are

$$\lambda_1 = 7, \quad \lambda_2 = 0$$

and the two normalized eigenvectors are

$$\mathbf{u}_1 = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}, \quad \mathbf{u}_2 = \begin{pmatrix} -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}.$$



28

The matrix that describes the K-L transformation is given by

$$U = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix} = \frac{1}{\sqrt{2}} \cdot \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} = \sqrt{2} \cdot \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{pmatrix}. \quad (3.124)$$

The K-L transformation maps the two dimensional data set Ψ in one dimension because λ_2 is zero (see Figure 3.35). For example, the data point $(1, 1)$ is mapped on the x -axis

$$\begin{pmatrix} \sqrt{2} \\ 0 \end{pmatrix} = \begin{pmatrix} \sqrt{2} \cdot \frac{1+1}{2} \\ \sqrt{2} \cdot \frac{1-1}{2} \end{pmatrix} = \sqrt{2} \cdot \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (3.125)$$

with the value $\sqrt{2} \approx 1.4142$ corresponding to the length of the vector $(1, 1)$.



29

Principal component analysis

- The variance in the direction of the k^{th} eigenvector (or principal component) is given by the eigenvalue λ_k
- Eigenvalues can be used to estimate how many components to keep
- **Rule of thumb:** keep enough to explain 85% of the variation
- $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$

$$\frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^n \lambda_j} \approx 0.85$$

if $k = n$, we preserve 100% of the original variation

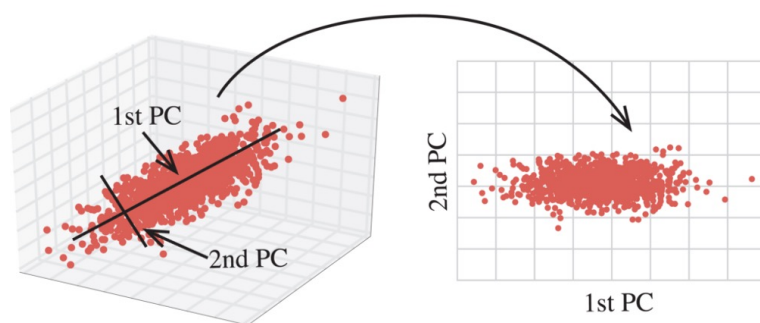
30

Problems

- Principal components are linear transformation of the original features
- It is difficult to attach any semantic meaning to principal components
- For new data which is added to the dataset, the PCA has to be recomputed

31

Space transformation



Axes of greater variance given by *eigenvectors* of covariance matrix

32

Reconstruction

- This PCA matrix W corresponds to a linear mapping from $\mathbb{R}^m \rightarrow \mathbb{R}^s$

$$\mathbf{z} = \mathbf{W}^T \cdot \mathbf{x}$$

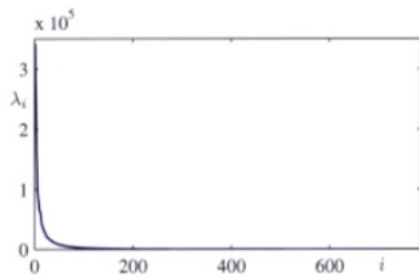
- Reconstruction

$$\mathbf{W} \cdot \mathbf{z} = \hat{\mathbf{x}}$$

- PCA minimizes the reconstruction error: $\|\mathbf{x} - \hat{\mathbf{x}}\|$

- It can be shown that the reconstruction error is: $error = 1/2 \sum_{i=k+1}^m \lambda_i$

33



- Eigenvalues



- Projection of the Eigenvectors, blue positive values, yellow negative values

PCA: Only the images of three

7 2 1 0 4 1 4 9 5 9
 0 6 9 0 1 5 9 7 8 4
 9 6 6 5 4 0 7 4 0 1
 3 1 3 4 7 2 7 1 2 1
 1 7 4 2 3 5 1 2 4 4
 6 3 5 5 6 0 4 1 9 5
 7 8 9 3 7 4 6 4 3 0
 7 0 2 9 1 7 3 2 9 7
 7 6 2 7 8 4 7 3 6 1
 3 6 9 3 1 4 1 7 6 9

34

Singular Value Decomposition

- We can use SVD to perform PCA
 - SVD is **more numerically stable** if the columns are close to collinear
 - Factorize a Covariance Matrix $A:=C$
 - Difference: compute the eigenvectors out of $C C^T = C^T C = C C$, use U as before....

Any matrix A can be factorized as

$$A = U \cdot S \cdot V^T$$

U is a orthogonal matrix with orthonormal eigenvectors from $A \cdot A^T$

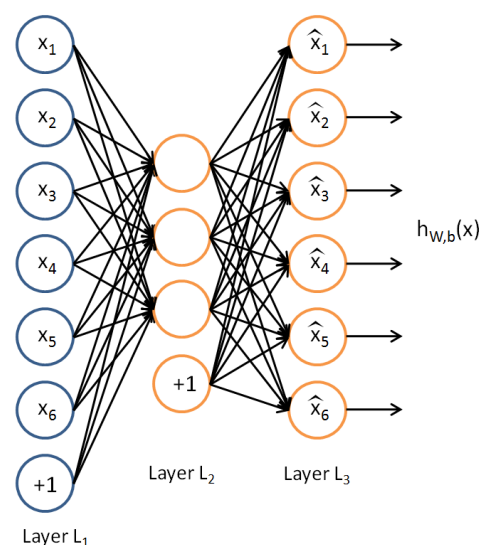
V a orthogonal matrix with orthonormal eigenvectors from $A^T \cdot A$

S is a diagonal matrix with r elements equal to the root of the positive eigenvalues of $A \cdot A^T$ or $A^T \cdot A$

35

Other Methods for Dimension Reduction

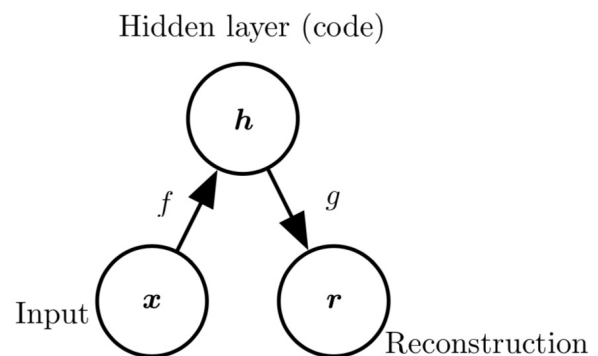
- Unsupervised Learning
 - Data: no labels!
 - Goal: Learn the structure of the data
- Traditionally, autoencoders were used for dimensionality reduction or feature learning.
- Undercomplete autoencoder (less units in the hidden layer)



36

Structure of an Autoencoder

- An autoencoder is a neural network that is trained to attempt to copy its input to its output.
- It has a hidden layer \mathbf{h} that describes a code used to represent the input.
- The network may be viewed as consisting of two parts:
 - An encoder function $\mathbf{h} = f(\mathbf{x})$
 - A decoder that produces a reconstruction $\mathbf{r} = g(\mathbf{h})$.



37

Avoiding Trivial Identity

- Undercomplete autoencoders
 - \mathbf{h} has lower dimension than \mathbf{x}
 - f or g has low capacity (e.g., linear g)
 - Must discard some information in \mathbf{h}
- Overcomplete autoencoders
 - \mathbf{h} has higher dimension than \mathbf{x}
 - Must be regularized

38

Undercomplete Autoencoders

- An autoencoder whose code dimension is less than the input dimension is called undercomplete
- Minimising the Loss function

$$L(\mathbf{x}, g(f(\mathbf{x})))$$

such as for example

$$E(\mathbf{w}) = \frac{1}{2} \cdot \sum_{k=1}^N (\mathbf{x}_k - g(f(\mathbf{x}_k)))^2$$

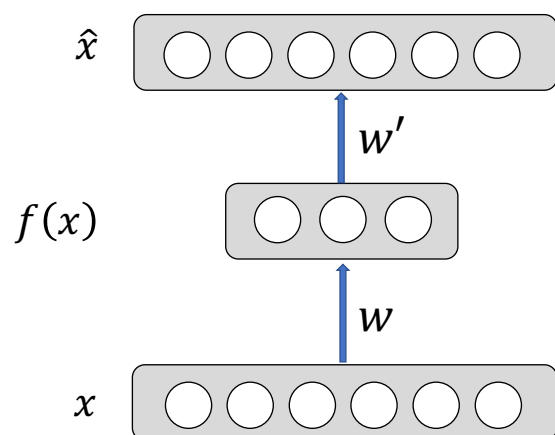
- If our input is interpreted as *bit vectors* or vectors of bit probabilities the *cross entropy* can be used

$$E(\mathbf{w}) = H(p, q) = \sum_{k=1}^N \sum_{bits} \mathbf{x}_k \cdot \log(g(f(\mathbf{x}_k)))$$

39

Undercomplete AE

- Hidden layer is **Undercomplete** if smaller than the input layer
 - Compresses the input
 - Compresses well only for the training dist.
- Hidden nodes will be
 - Good features for the training distribution.
 - Bad for other types on input

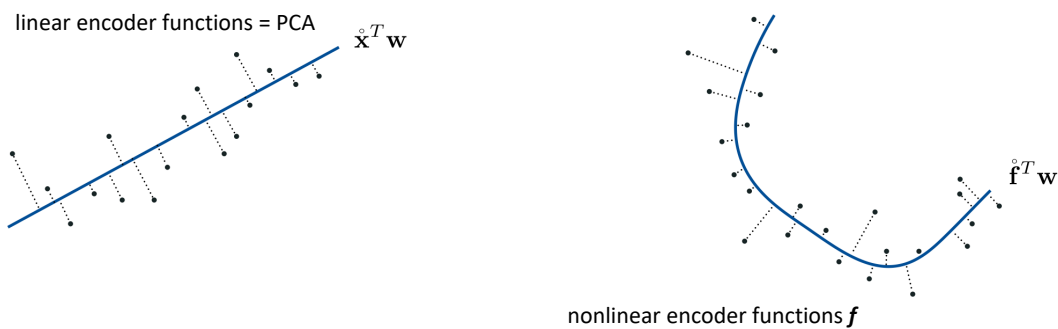


40

$$L(\mathbf{x}, g(f(\mathbf{x})))$$

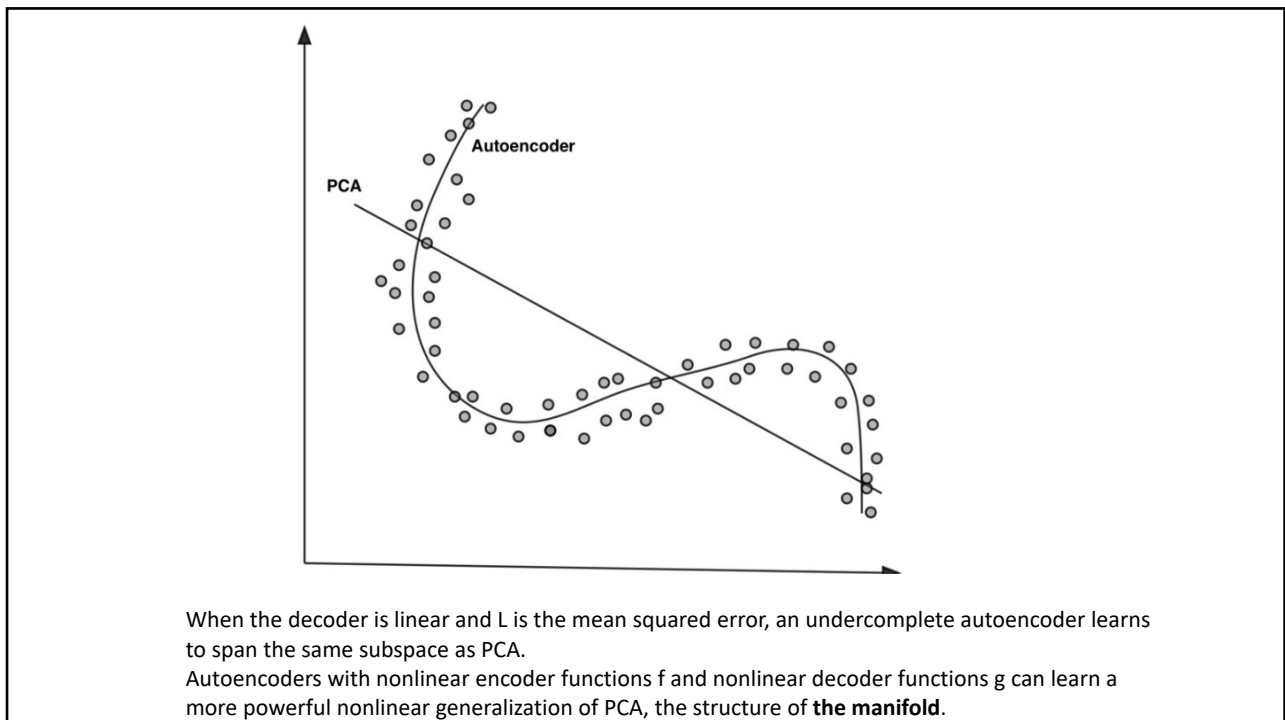
- trained with backpropagation using minibatches
- Learning an undercomplete representation forces the autoencoder to capture the **most important features** of the training data.
- When the decoder is linear ($f(x)$) and Loss is the mean squared error, an under- complete autoencoder learns to span the same subspace as PCA

41



- Autoencoders with nonlinear encoder functions $\mathbf{f}(\mathbf{x})$ and nonlinear decoder functions \mathbf{g} can thus learn a more powerful nonlinear generalization of PCA

42



43

Too much Capacity

- An autoencoder with a one-dimensional code but a very powerful **nonlinear encoder** could learn to represent each training example $\mathbf{x}^{(i)}$ with the **code i** .
- The decoder could learn to map these integer indices back to the values of specific training examples.
- An autoencoder trained to perform the copying task can fail to learn anything *useful* about the dataset if the **capacity** of the autoencoder is allowed to become **too great**.

44

Overcomplete Autoencoders

- An autoencoder whose code dimension is bigger than the input dimension is called overcomplete
- For Overcomplete Autoencoders a linear encoder and linear decoder can learn to copy the input to the output without learning anything useful about the data distribution.
- A **regularized** autoencoder can be nonlinear and overcomplete but still learn something useful about the data distribution even if the model capacity is great enough to learn a trivial identity function.
- Leads to *sparse autoencoder*.

45

Sparse Autoencoders

- Limit capacity of autoencoder by adding a term to the cost function penalizing the code for being larger
- Special case of variational autoencoder
- Probabilistic model
- Laplace prior corresponds to l_1 sparsity penalty Dirac variational posterior

46

Sparse Autoencoders

$$L(\mathbf{x}, g(f(\mathbf{x})) + \Omega(\mathbf{h}))$$

with sparsity penalty $\Omega(\mathbf{h})$, $g(\mathbf{h})$ is the decoder output and typically we have

$$\mathbf{h} = f(\mathbf{x})$$

For example

$$E(\mathbf{w}) = \frac{1}{2} \cdot \sum_{k=1}^N (\mathbf{x}_k - g(f(\mathbf{x}))_k)^2 + \lambda \cdot \|\mathbf{w}\|_1$$

for one hidden layer. We can use $\lambda \cdot \|\mathbf{w}\|_1$ only for \mathbf{h}

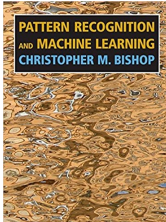
47

Sparse Autoencoders

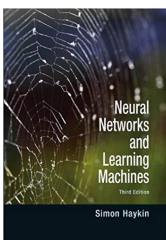
- An autoencoder that has been regularized to be sparse must respond to unique statistical features of the dataset it has been trained on, rather than simply acting as an identity function.
- In this way, training to perform the copying task with a **sparsity penalty** can yield a model that has learned useful features
- One way to achieve actual zeros in \mathbf{h} for sparse autoencoders is to use rectified linear units (*ReLU*) to produce the code layer.
- With a prior that actually pushes the representations to zero $\lambda \cdot \|\mathbf{w}\|_1$, one can indirectly control the average number of zeros in the representation.

48

Literature



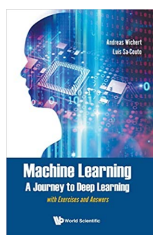
- Christopher M. Bishop, Pattern Recognition and Machine Learning (Information Science and Statistics), Springer 2006
 - Chapter 12



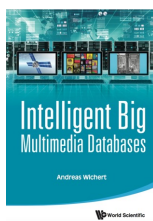
- Simon O. Haykin, Neural Networks and Learning Machine, (3rd Edition), Pearson 2008
 - Chapter 10

49

Literature



- Machine Learning - A Journey to Deep Learning, A. Wichert, Luis Sa-Couto, World Scientific, 2021
 - Chapter 15



- Intelligent Big Multimedia Databases, A. Wichert, World Scientific, 2015
 - Chapter 3, Section 3.3

50