

A tutorial on principal component analysis

Rasmus R. Paulsen

DTU Compute

Based on

Jonathan Shlens: A tutorial on Principal Component Analysis (version 3.02 – April 7, 2014)

http://compute.dtu.dk/courses/02503



2025

What is your experience with Principal Component Analysis (PCA)



I never heard of PCA before this course

I have seen PCA mentioned before

I have read about PCA but never used it

I have used PCA a few times

PCA and I are practically best friends

What is your experience with Principal Component Analysis (PCA)



0%

I never heard of PCA before this course 0%

I have seen PCA mentioned before

I have read about PCA but never used it 0%

I have used PCA a few times 0%

PCA and I are practically best friends 0%

What is your experience with Principal Component Analysis (PCA)



I never heard of PCA before this course 0%

I have seen PCA mentioned before

I have read about PCA but never used it

I have used PCA a few times

PCA and I are practically best friends

0%

0%

0%

0%



Principal Component Analysis (PCA) learning objectives

- Describe the concept of principal component analysis
- Explain why principal component analysis can be beneficial when there is high data redundancy
- Arrange a set of multivariate measurements into a matrix that is suitable for PCA analysis
- Compute the covariance of two sets of measurements
- Compute the covariance matrix from a set of multivariate measurements
- Compute the principal components of a data set using Eigenvector decomposition
- Describe how much of the total variation in the data set that is explained by each principal component



Iris data

The Iris flower data set or Fisher's Iris data set is a data set introduced by Ronald Fisher in his 1936 paper The use of multiple measurements in taxonomic problems











Iris data



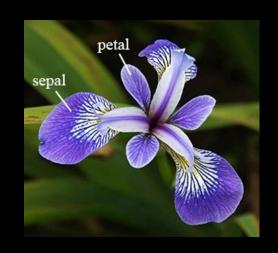
- 3 Iris types
 - 50 flowers of each type
- For each flower
 - Sepal length
 - Sepal width
 - Petal length
 - Petal width
- We use one type as example
 - 50 measured flowers

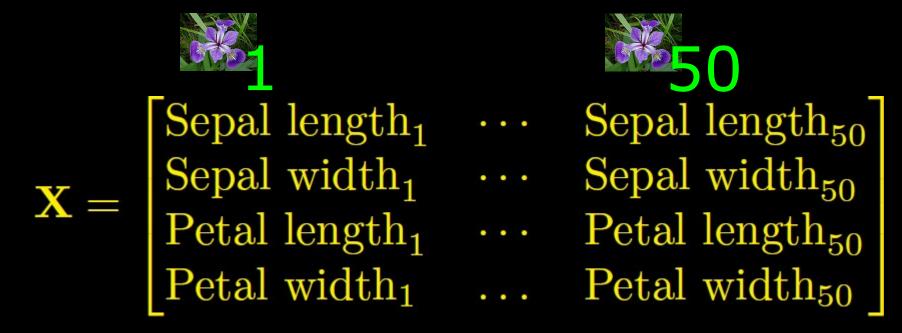




Iris Data Matrix

- One column is one flower
- One row is all measurements of one type









What can we use these data for?



- The measurements can be used to:
 - Recognize a species of flowers
 - Classify flowers into groups
 - Describe the characteristics of the flower
 - Quantify growth rates
 - ...
- Do we need all the measurements?
 - Can we boil down or combine some measurements?
- Are some measurements redundant?



2025



Variance

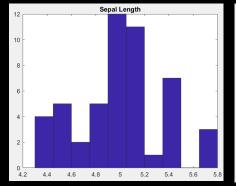
$$\sigma_{SL}^2 = 0.1242$$

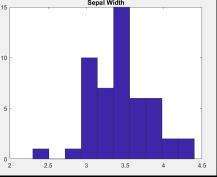
$$\sigma_{SW}^2 = 0.1437$$

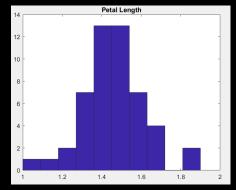
$$\sigma_{PL}^2 = 0.0302$$

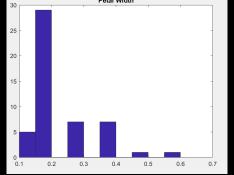
$$\sigma_{PW}^2 = 0.0111$$







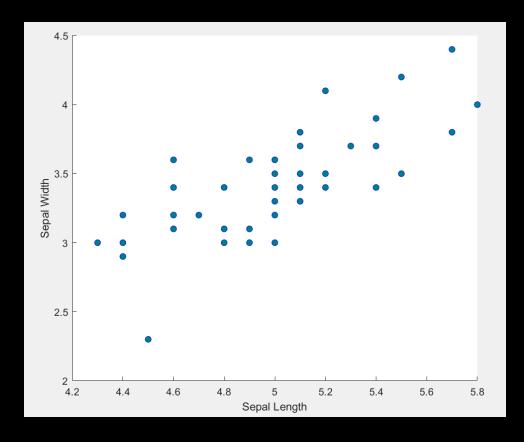




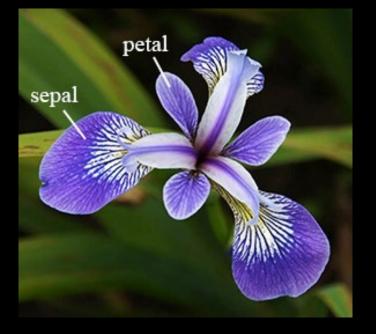




High Redundancy



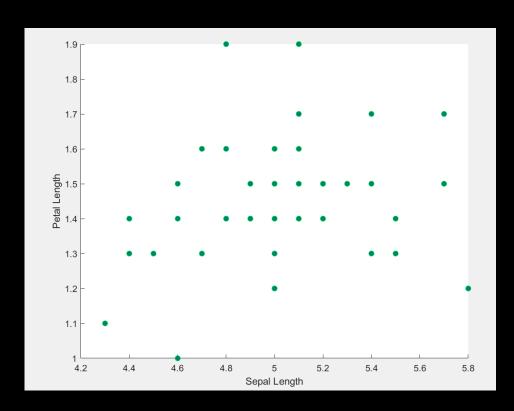
Observation: We can explain quite a lot of the sepal width if we know the sepal lengths



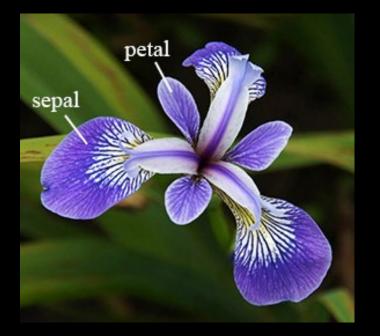




Low Redundancy



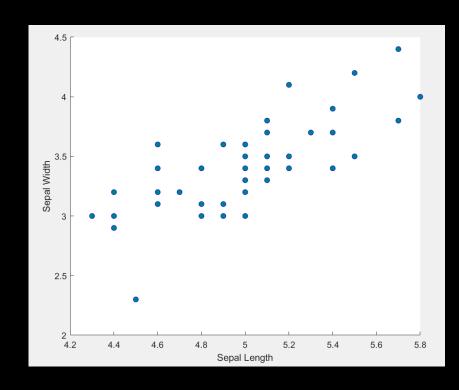
Observation: We can **NOT** explain the petal length if we know the sepal lengths



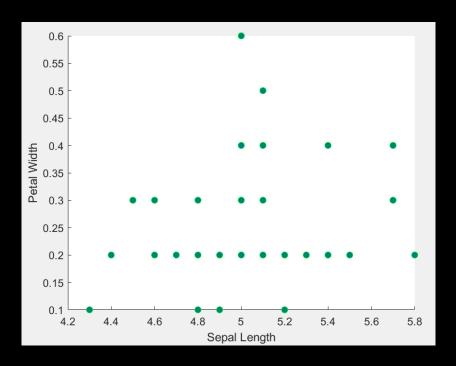




Covariance





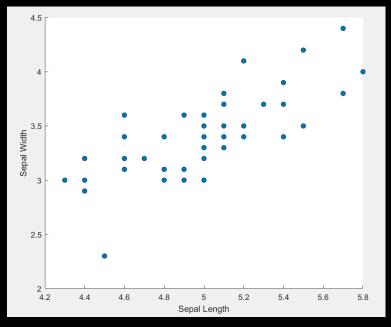


Covariance measures the *relationship* between measurements



÷

High Covariance





Sepal length and sepal width

$$a_i = SL = \{5.1, 4.9 ..., 5\}$$

$$b_i = SW = \{3.5, 3, ..., 3.3\}$$

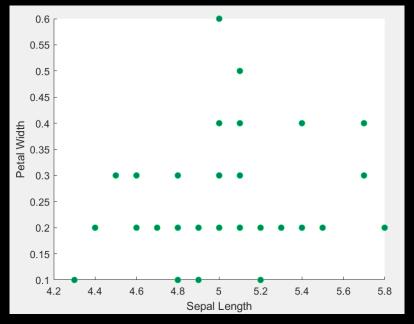
$$\sigma_{\text{SL,SW}}^2 = \frac{1}{n} \sum_{i} a_i b_i = 17.2578$$

Note that in practice n-1 is used instead of n





Low covariance





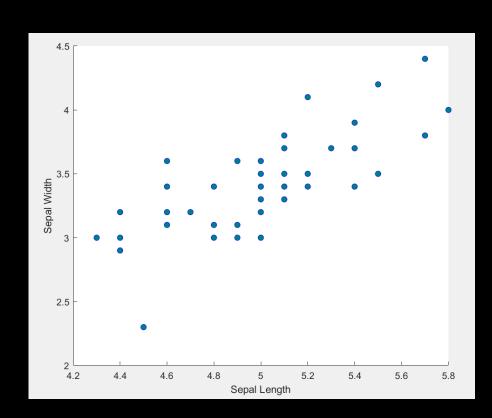
Sepal length and petal width

$$\sigma_{\rm SL,PW}^2 = \frac{1}{n} \sum_i a_i b_i = 1.2416$$





Vector notation for covariance



Sepal length and sepal width

$$a = SL = [5.1, 4.9 ..., 5]$$

$$\mathbf{b} = SW = [3.5, 3, ..., 3.3]$$

$$\sigma_{ ext{SL,SW}}^2 \;\; = \;\; rac{1}{n} \mathbf{a} \mathbf{b}^T$$





Matrix notation for covariance

$$m \times n$$
 matrix (m=4 and n=50)

$$\mathbf{X} = \begin{bmatrix} \text{Sepal length}_1 & \cdots & \text{Sepal length}_{50} \\ \text{Sepal width}_1 & \cdots & \text{Sepal width}_{50} \\ \text{Petal length}_1 & \cdots & \text{Petal length}_{50} \\ \text{Petal width}_1 & \cdots & \text{Petal width}_{50} \end{bmatrix}$$

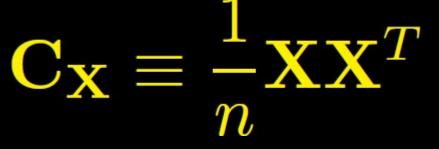
$$\mathbf{C}_{\mathbf{X}} \equiv \frac{1}{n} \mathbf{X} \mathbf{X}^T$$
 $m \times m \text{ square matrix } (m=4)$

Note that in practice n-1 is used instead of n



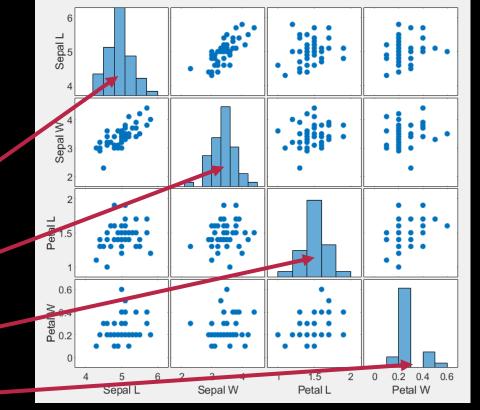


Covariance matrix autopsy



The diagonal elements are the variances

$$\sigma_{SL}^2 = 0.1242$$
 $\sigma_{SW}^2 = 0.1437$
 $\sigma_{PL}^2 = 0.0302$
 $\sigma_{PW}^2 = 0.0111$







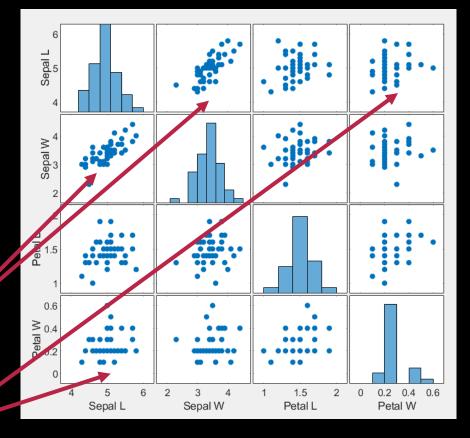
Covariance matrix autopsy II

$$\mathbf{C}_{\mathbf{X}} \equiv \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

The off-diagonal elements are the covariance

$$\sigma_{{
m SL,SW}}^2 = \frac{1}{n} \sum_i a_i b_i = 17.2578$$

$$\sigma_{\text{SL,PW}}^2 = \frac{1}{n} \sum_{i} a_i b_i = 1.2416$$



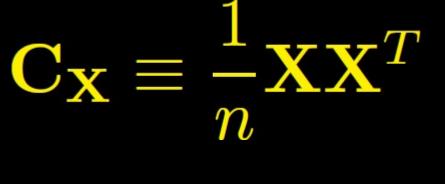
Symmetric!



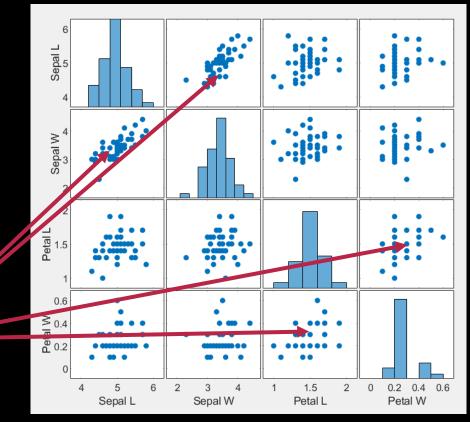
2025



Covariance matrix autopsy III



High redundancy



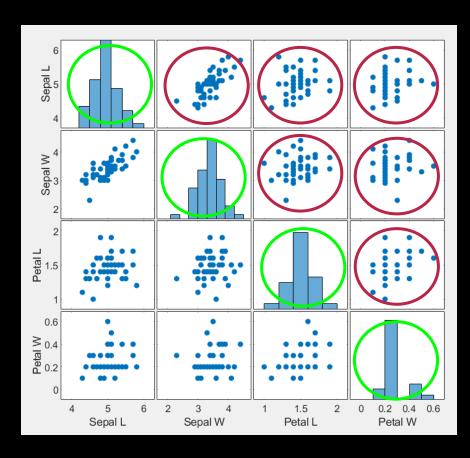
Symmetric!





Goals

21



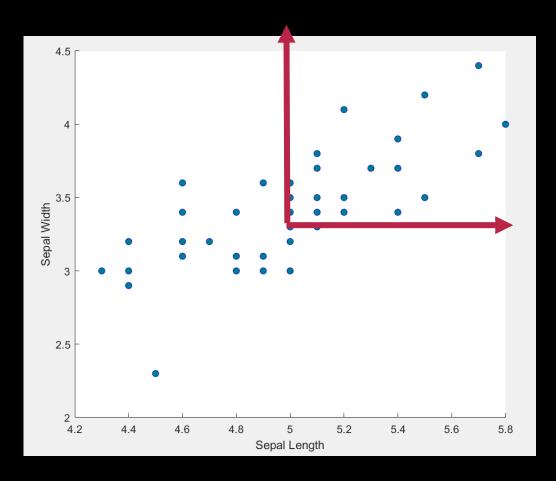
- Minimize redundancy
 - Covariance should be small
- Maximize signal
 - Variance should be large

Signal to noise ratio:

$$SNR = \frac{\sigma_{signal}^2}{\sigma_{noise}^2}$$



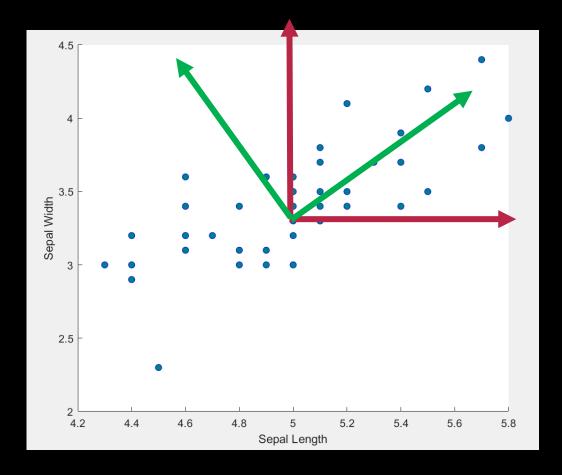




- We start by subtracting the mean
 - Centering data
- Red lines are the default basis

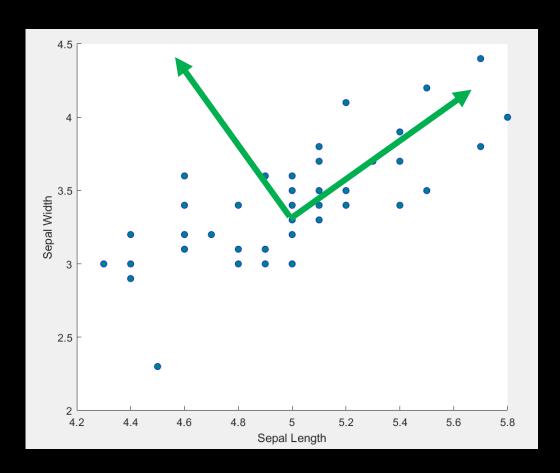








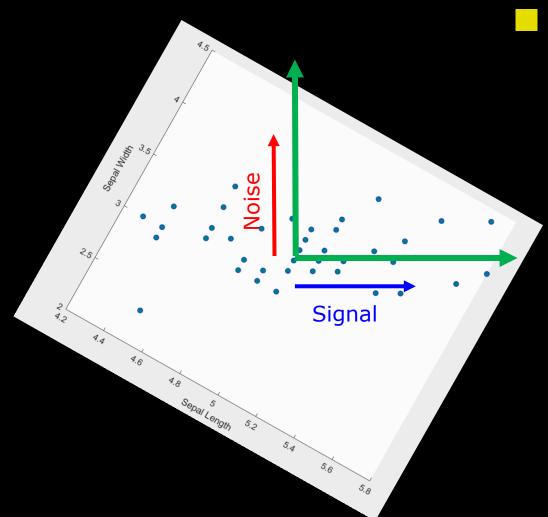




A new basis that follows the covariance in the data



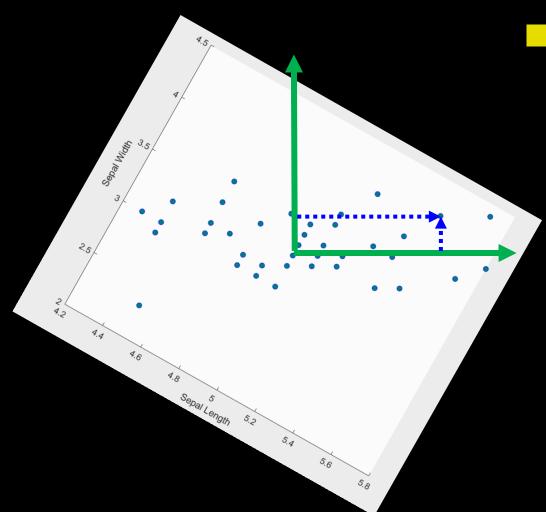




Lets try to rotate the data – for visualisation



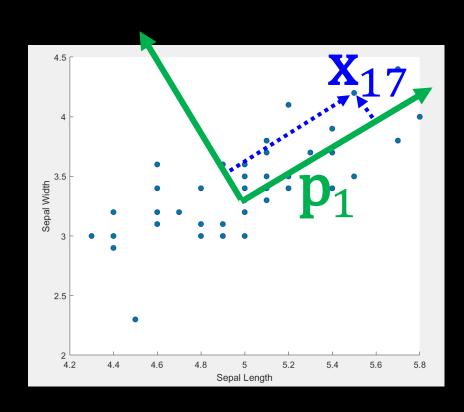




Finding the measurement values in the new basis





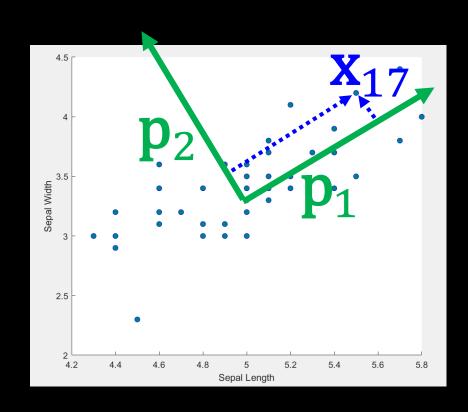


The dot product projects a point down to a new axis

$$\mathbf{x}_{17,\text{new}} = x_{17} \cdot p_1$$







The dot product projects a point down to a new axis

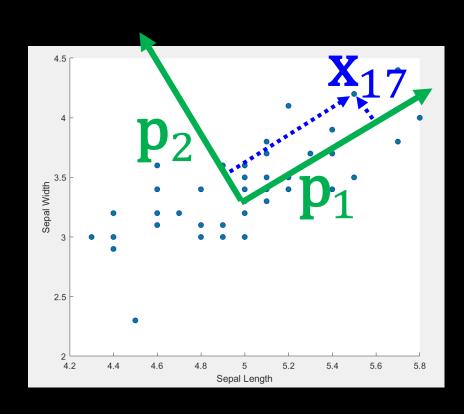
$$\mathbf{PX} = \mathbf{Y}$$

 p_1 p₁ and p_2 are the rows of P

$$\mathbf{X} = \begin{bmatrix} \text{Sepal length}_1 & \cdots & \text{Sepal length}_{50} \\ \text{Sepal width}_1 & \cdots & \text{Sepal width}_{50} \\ \text{Petal length}_1 & \cdots & \text{Petal length}_{50} \\ \text{Petal width}_1 & \cdots & \text{Petal width}_{50} \end{bmatrix}$$







The dot product projects a point down to a new axis

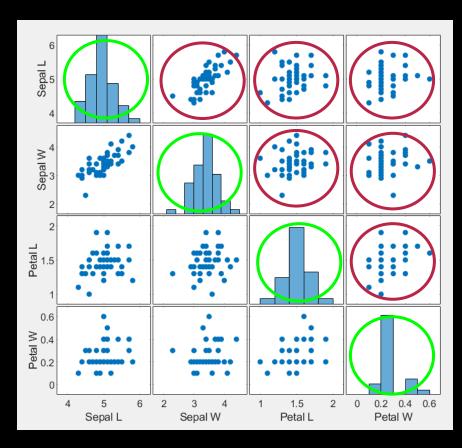
$$\mathbf{PX} = \mathbf{Y}$$

Here Y contains the new coordinates/measurements per sample





Goals



- Minimize redundancy
 - Covariance should be small
- Maximize signal
 - Variance should be large
- Transform our data
 - Rotating and scaling the basis

$$\mathbf{Y} = \mathbf{P}\mathbf{X}$$

So it will have

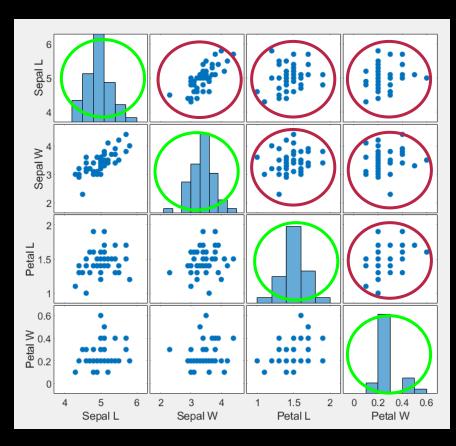
$$\mathbf{C}_{\mathbf{Y}} \equiv \frac{1}{n} \mathbf{Y} \mathbf{Y}^T$$





Goals

31



■ The covariance matrix

$$\mathbf{C}_{\mathbf{Y}} \equiv \frac{1}{n} \mathbf{Y} \mathbf{Y}^T$$

- Should be as diagonal as possible
- We do this by

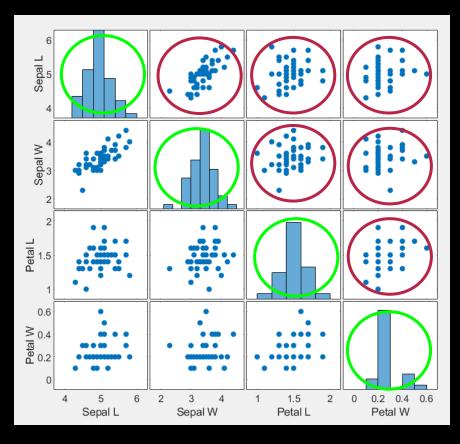
$$Y = PX$$

Where **P** are the principal components





Computing the principal components



The Principal Components of X are the eigenvectors of

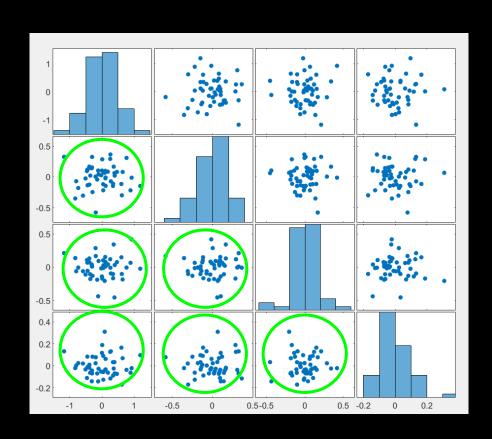
$$\mathbf{C}_{\mathbf{X}} \equiv \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

The i'th diagonal value of C_V is the variance along principal component number i





New covariance matrix for Iris data



The principal component are found and

$$Y = PX$$

With the covariance matrix

$$\mathbf{C}_{\mathbf{Y}} \equiv \frac{1}{n} \mathbf{Y} \mathbf{Y}^T$$

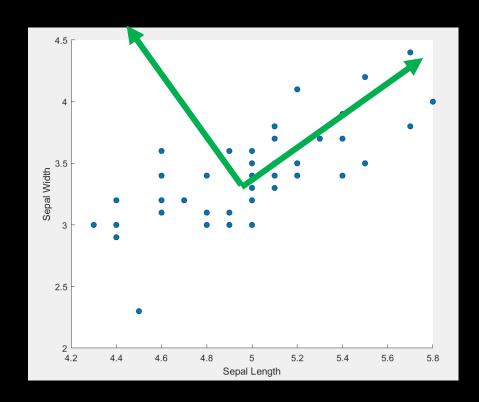
Covariance: 0

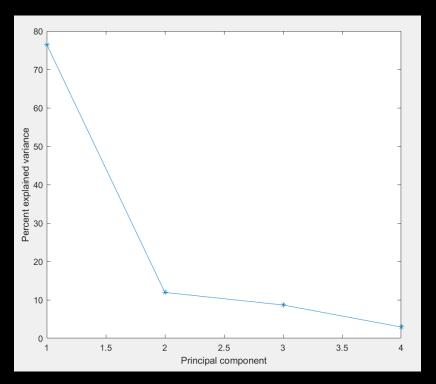


2025



Explained variance





One component explains 75% of the total variation – so for each flower we can have one number that explains 75% percent of the 4 measurements!





What can we use it for?

Classification



Based on one value instead of 4





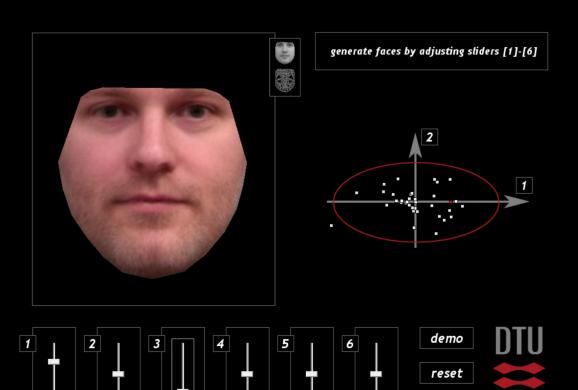


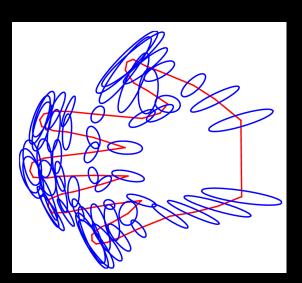


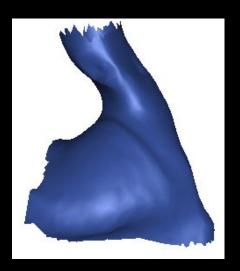
What can we use it for?

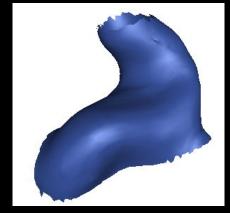
Many more examples in the course

help













Final note – practical estimation of covariance matrix

$$\mathbf{C}_{\mathbf{X}} \equiv \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

In practice n-1 is used instead of n for exercises and in the exam.

$$\mathbf{C}_{\mathbf{X}} \equiv \frac{1}{n-1} \mathbf{X} \mathbf{X}^T$$

