

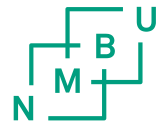
INF250

Multivariate analysis



Topics for INF250

- Multivariate analysis
 - PCA
 - Clustering
-



The world is multivariate

Multivariate analysis refers to any statistical technique used to analyze data that arises from more than one variable



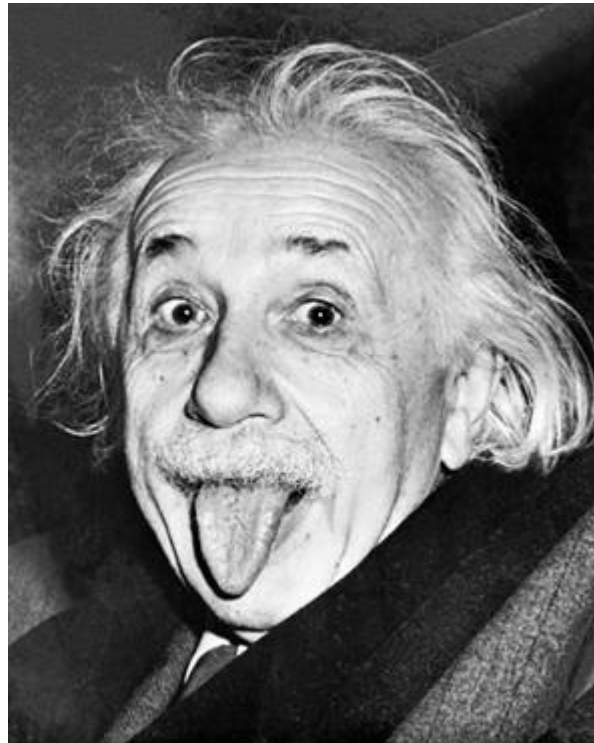
Who is this ?



WHO IS THIS?



WHO IS THIS

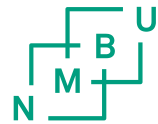


- The brain do not collect information from single parts, but for elements in a pattern
 - We need multivariate techniques!
-

- The diagram shows two data matrices, **X** and **Y**, representing the input and output data for a neural network.

Matrix X (X-variables): This matrix has n rows (objects) and p columns (variables). The elements are denoted as X_{ij} , where i is the object index and j is the variable index. The matrix is labeled **X** in the center.

Matrix Y (Y-variables): This matrix has n rows (objects) and c columns (variables). The elements are denoted as Y_{ij} , where i is the object index and j is the variable index. The matrix is labeled **Y** in the center.



Data presentation and modelling

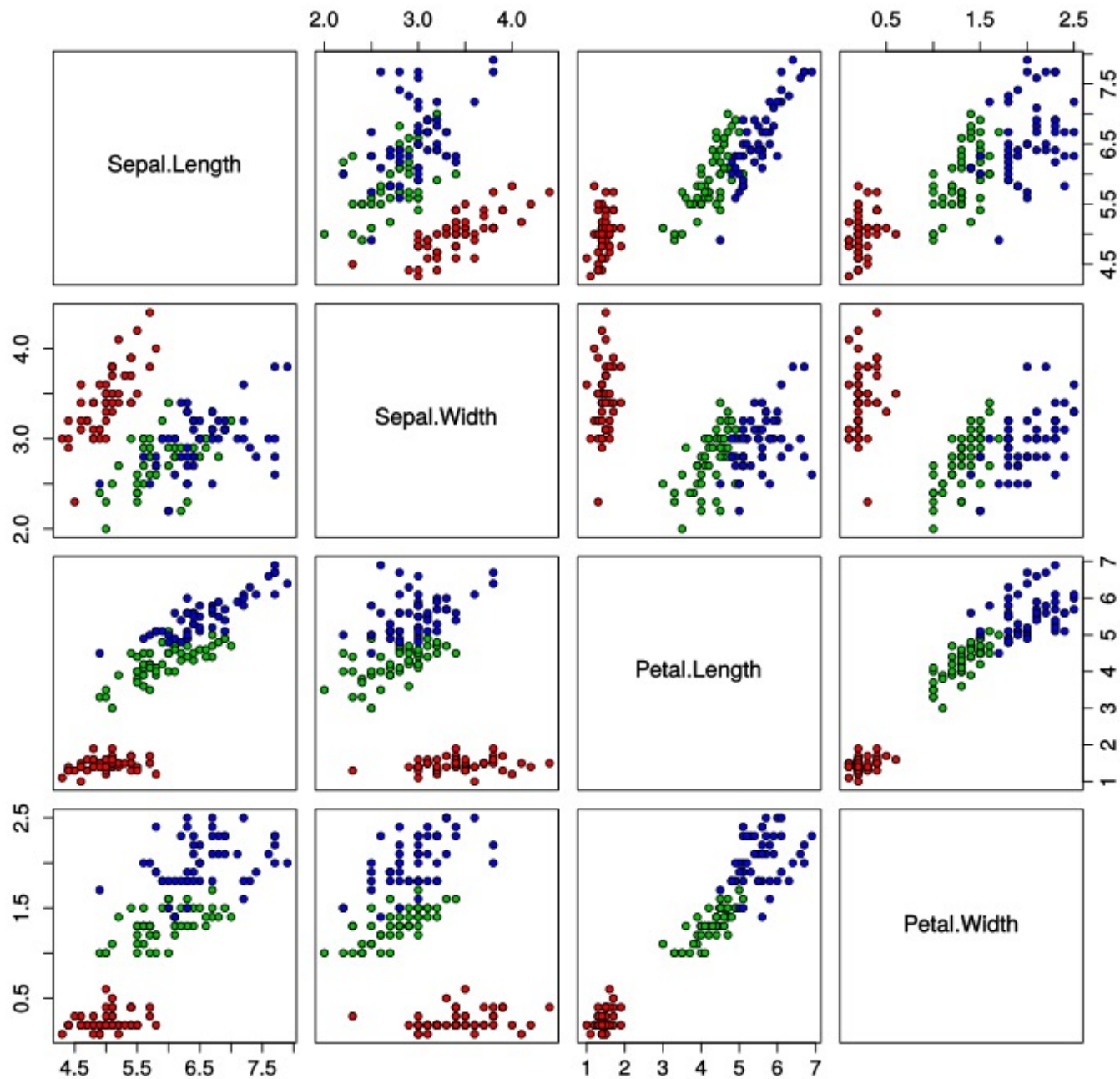
- One, Two, Three dimensional plots?
 - Do we need a N dimension space to view data?
 - How to find the optimal low dimension space that reveals maximum useful information?
 - How to find the «not so useful information»
 - «NOISE»
-



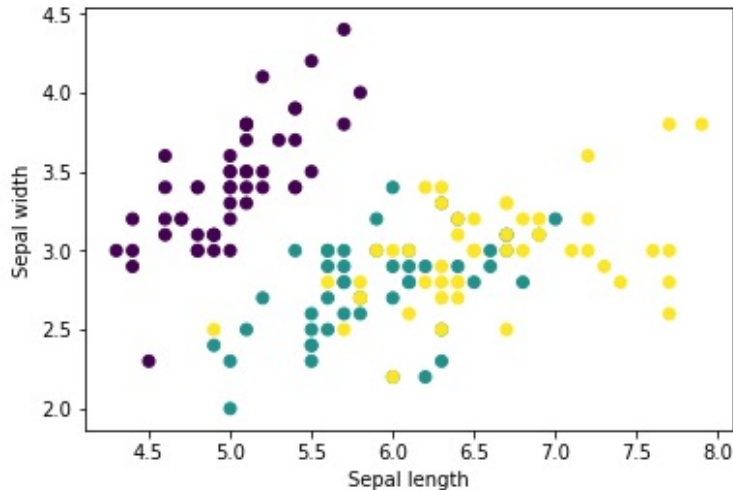
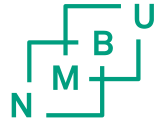
The Iris data set

- [https://en.wikipedia.org/wiki/Iris flower data set](https://en.wikipedia.org/wiki/Iris_flower_data_set)
 - Often used as a demo dataset for multivariate analysis
 - 3 species of Iris determined from 4 features: length and width of sepal and petals
 - http://scikit-learn.org/stable/auto_examples/datasets/plot_iris_dataset.html
-

Iris Data (red=setosa,green=versicolor,blue=virginica)



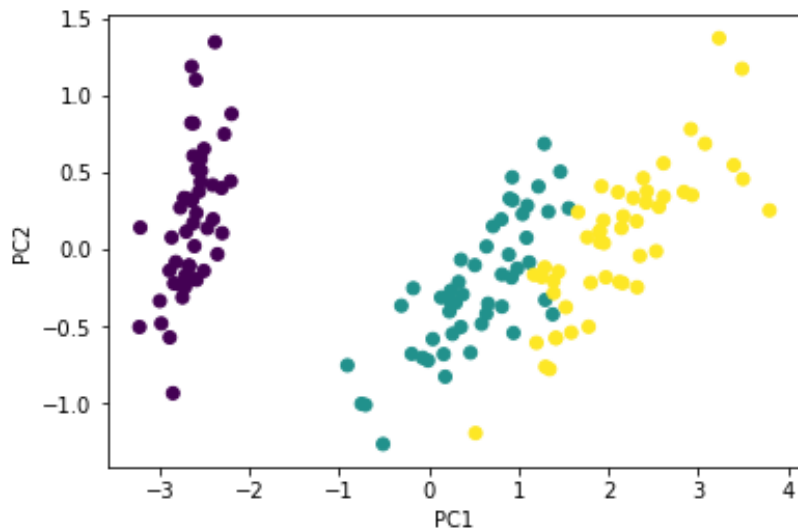
Principal Component Analysis



```
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.decomposition import PCA
```

```
# import the Iris data
iris = datasets.load_iris()
X = iris.data
y = iris.target
```

```
plt.figure()
plt.scatter(X[:,0],X[:,1],c=y)
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')
plt.show()
```

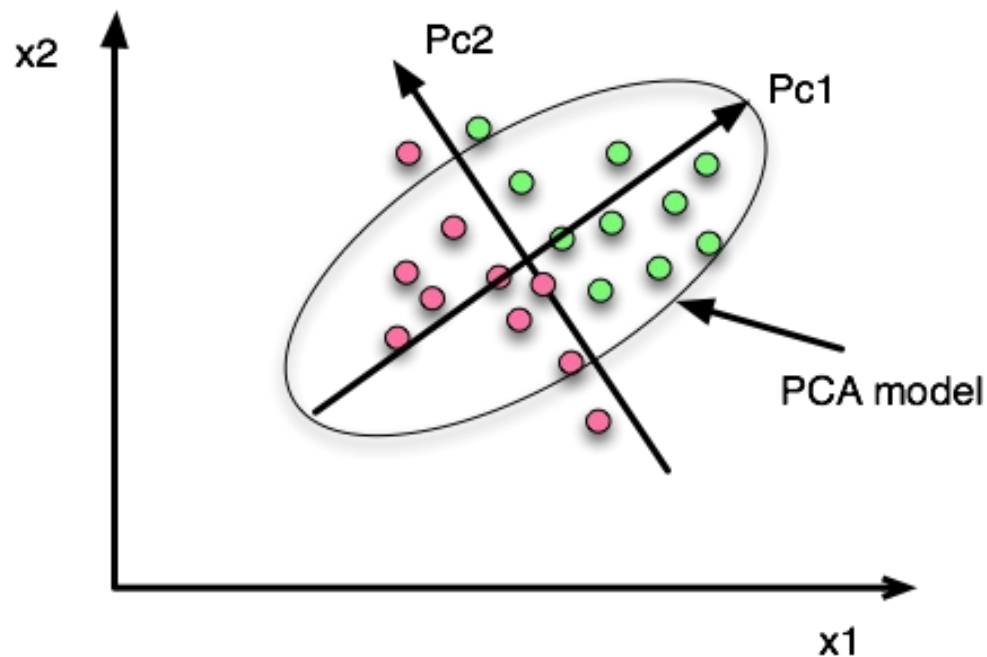


```
#simple PCA
```

```
X_reduced = PCA(n_components=3).fit_transform(iris.data)
plt.figure()
plt.scatter(X_reduced[:,1],X_reduced[:,2],c=y)
plt.xlabel('PC1')
plt.ylabel('PC2')
plt.show()
```

GRAPHICAL INTERPRETATION OF PCA

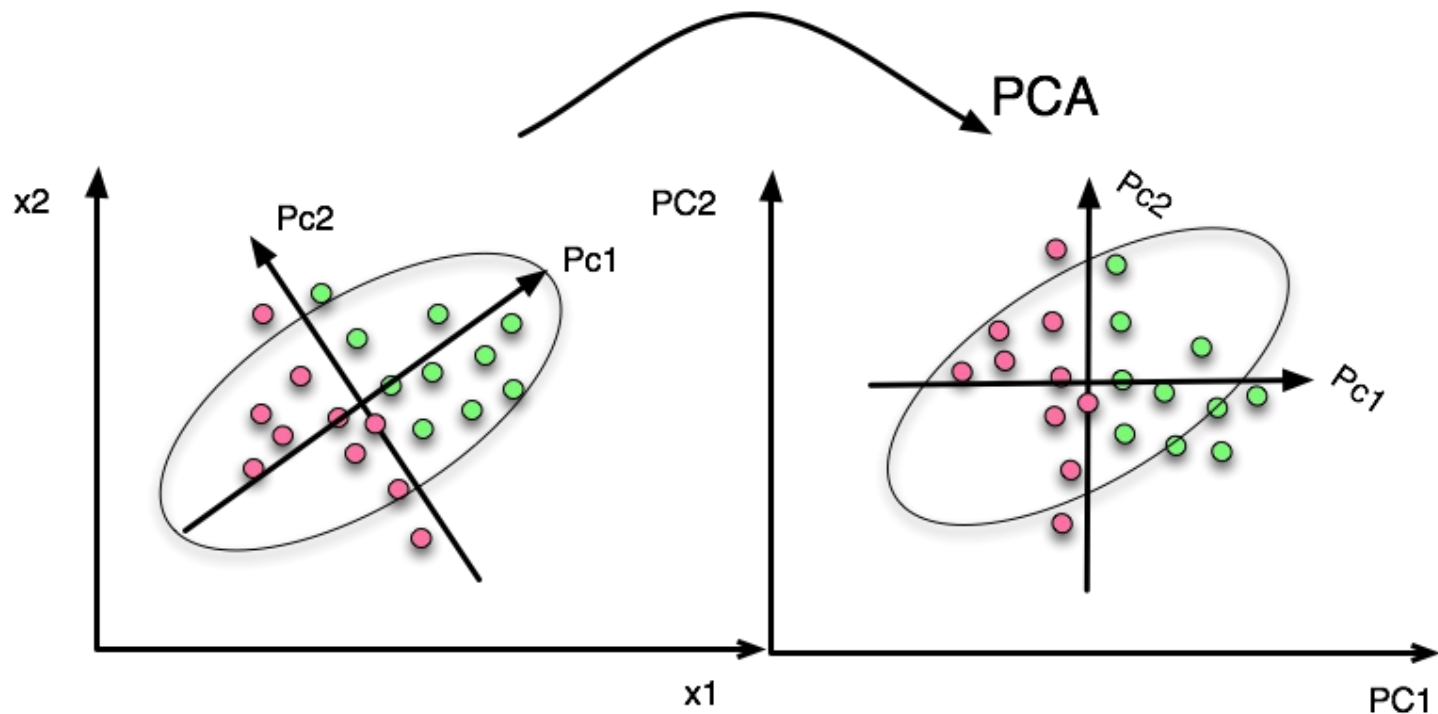
- The main idea is to form a minimum number of new variables that describes the variation of the data by a linear combination of the original values



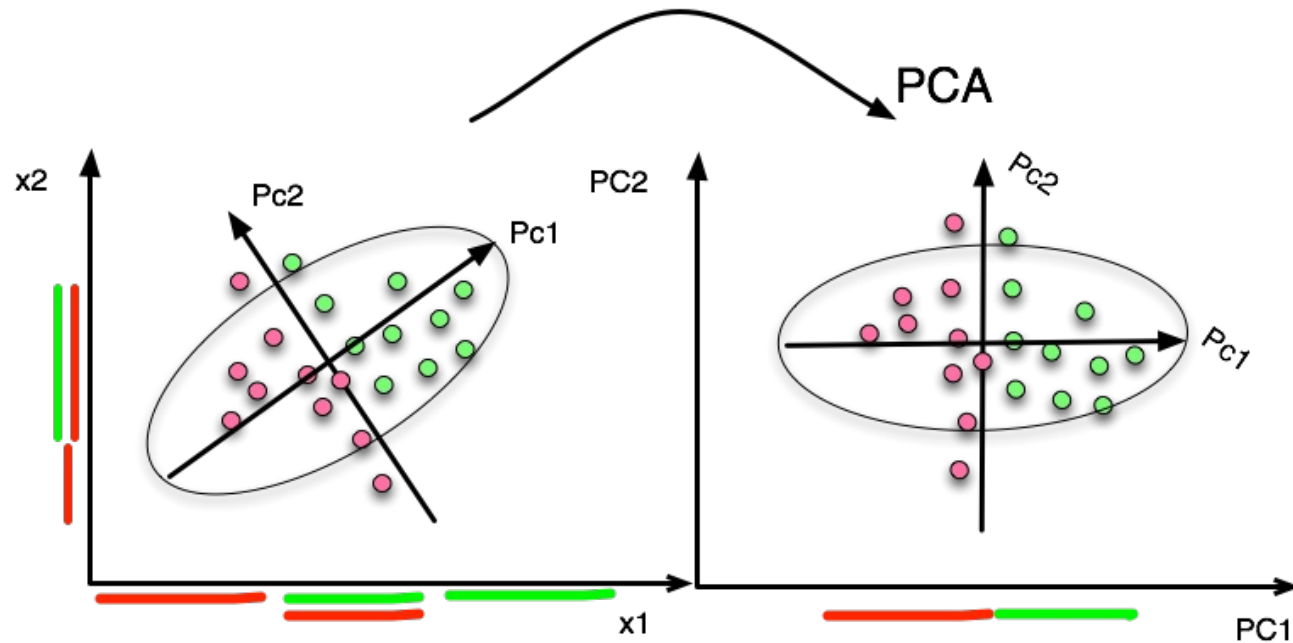
- PC₁
 - Direction of the largest variation in X
- PC₂
 - Direction of the next largest variation in X

The Principal Components

- The new coordinate system represents a stepwise rotation along the principal orthogonal axes



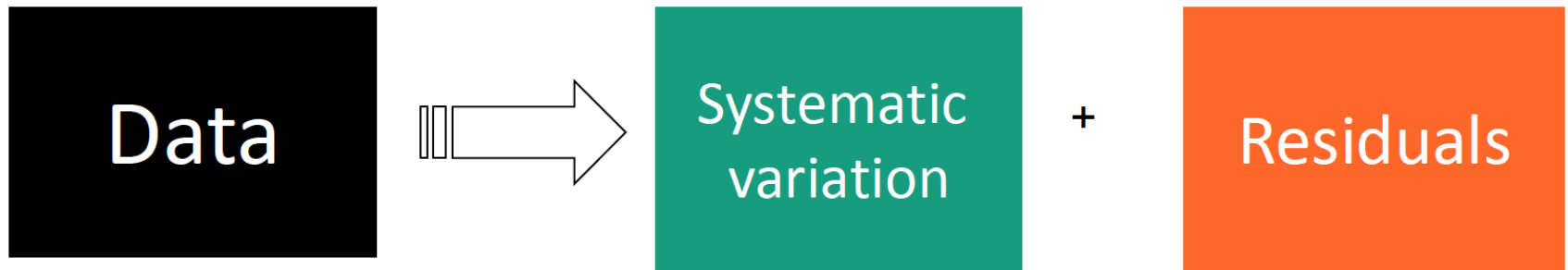
USE PCA TO SEPARATE CLASSES



Saparation of classes with a rotation

PCA basics

$$\begin{matrix} & & U \\ & & | \\ \text{---} & \text{---} & R \\ & & | \\ & & U \\ & & | \\ & & U \\ & & | \\ N & \text{---} & M & \text{---} & B & \text{---} & U \end{matrix}$$



$$X = TP' + E$$



X : data matrix

T : PCA scores matrix

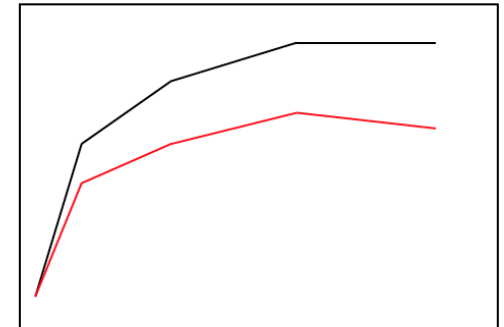
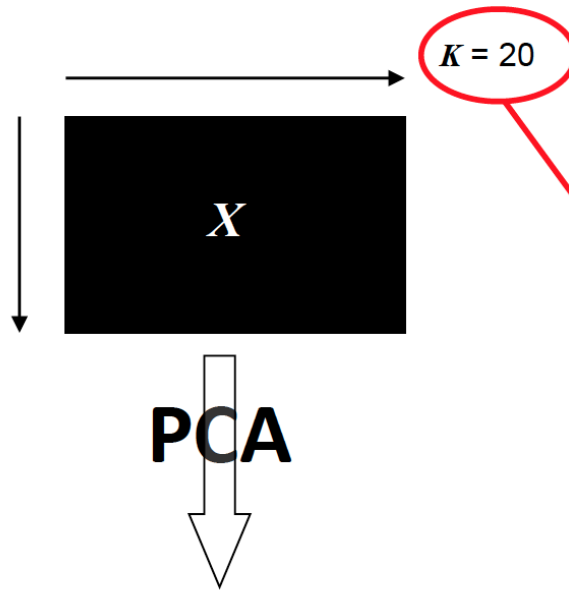
P : PCA loadings matrix

E : residuals / noise

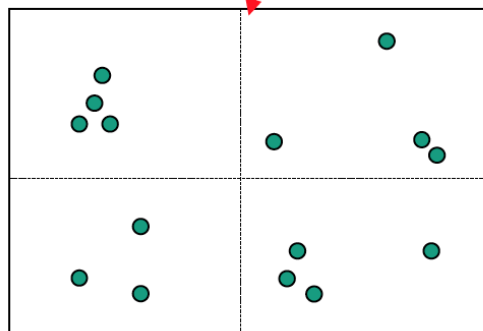
Principal Components (PC's) describing
the systematic variation in the data

PCA basics

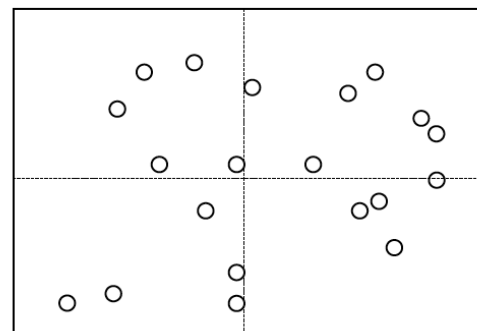
Data in this illustration consists of 15 observations and 20 variables → data matrix X of dimension (15×20)



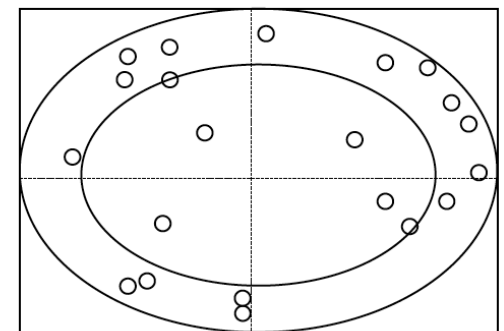
Explained variance plot



Scores plot



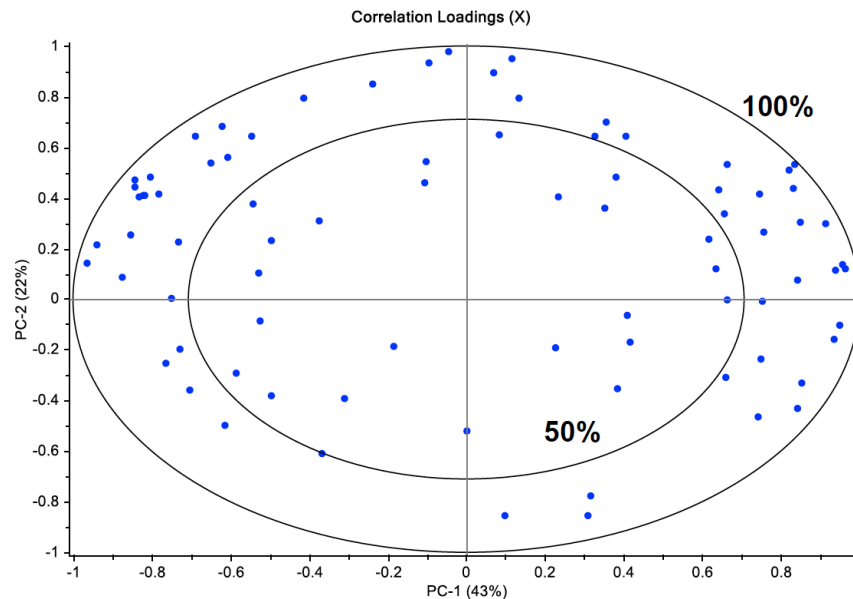
Loadings plot



Correlation loadings plot

PCA- correlation loadings

- Circles in the plot corresponding to various degrees of explained variances
- Typically one will present a circle for **100%** explained and for **50%** explained variance by the **two** components



PCA basics – centring and standardisation



- Only purpose of PCA is to look for directions with **high variance**
 - This implies: if there are variables x_k in X that have a **larger variance** than others ...
 - they will be given **most** attention
 - → They will **dominate** the extracted components
 - → They will **dominate** the plots
 - Generally one is interested in letting all variables play a role in the estimation of components (there are exceptions) → standardise variables x_k in X
 - Matrices in multivariate statistics are always **either** *centered* or *standardised*
-

PCA basics – centring and standardisation



$$X = \begin{pmatrix} x_{11} & \cdots & x_{1K} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{NK} \end{pmatrix}$$

- Number of **objects (rows)**:
 - $n = 1 \dots N$
- Number of **variables (columns)**:
 - $k = 1 \dots K$
- Observed value x_{nk} for
 - n 'th object
 - k 'th variable

center

$$x_{nk,cent} = x_{nk} - \bar{x}_k$$

$$\bar{x}_k = \frac{1}{N} \sum_{n=1}^N x_{nk}$$

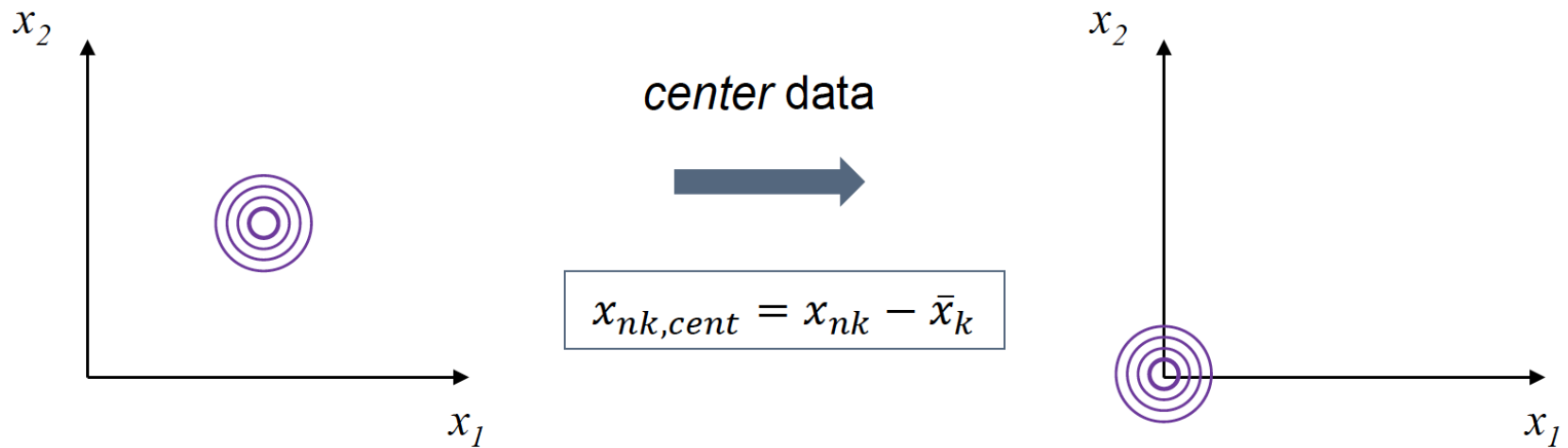
where

standardise

$$x_{nk,stand} = \frac{x_{nk} - \bar{x}_k}{\sigma_k}$$

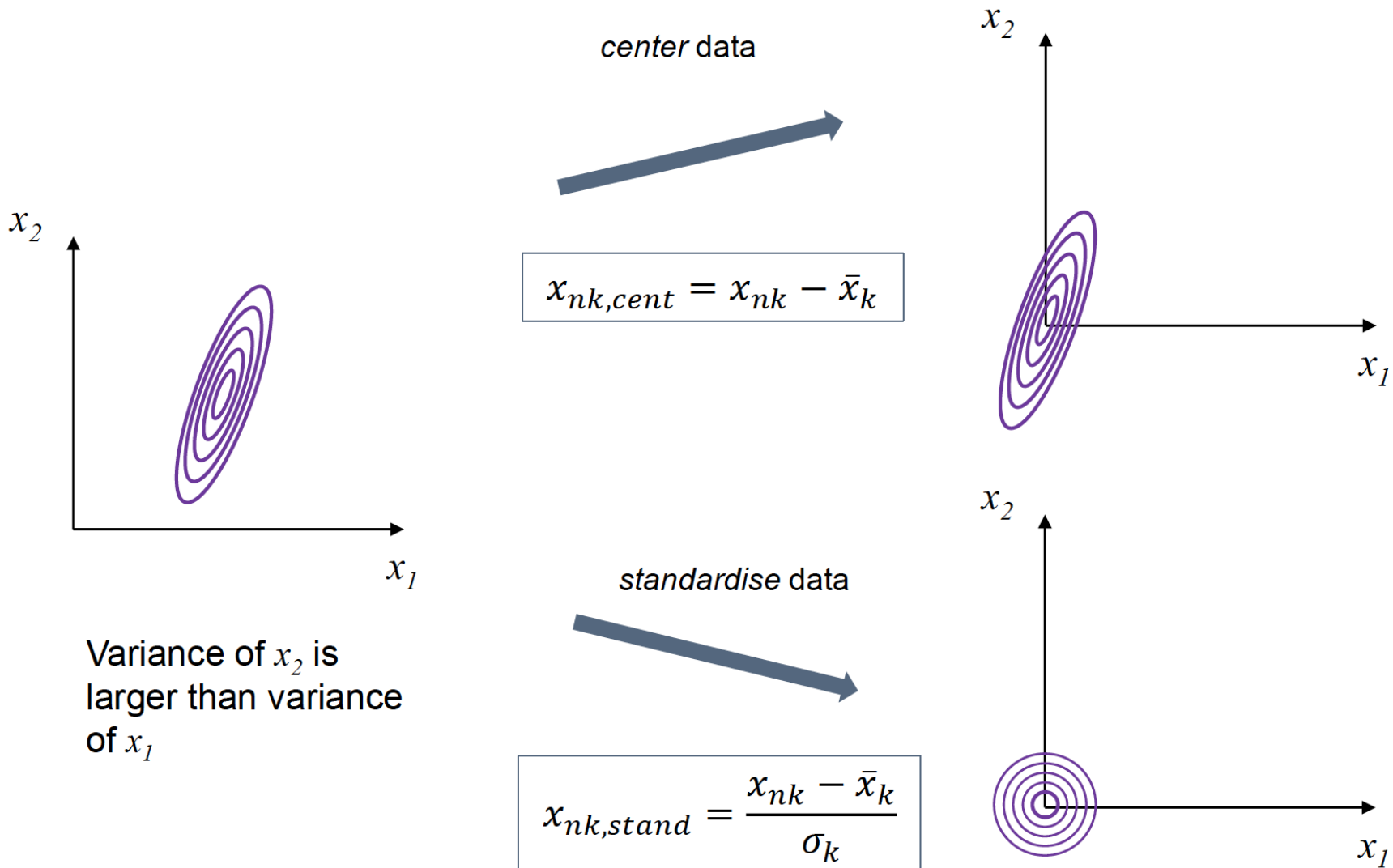
$$\sigma_k = \sqrt{\frac{1}{N} \sum_{n=1}^N (x_{nk} - \bar{x}_k)^2}$$

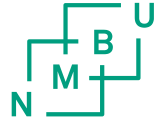
PCA basics – centring and standardisation



Equal variance of
 x_1 and x_2

PCA basics – centring and standardisation



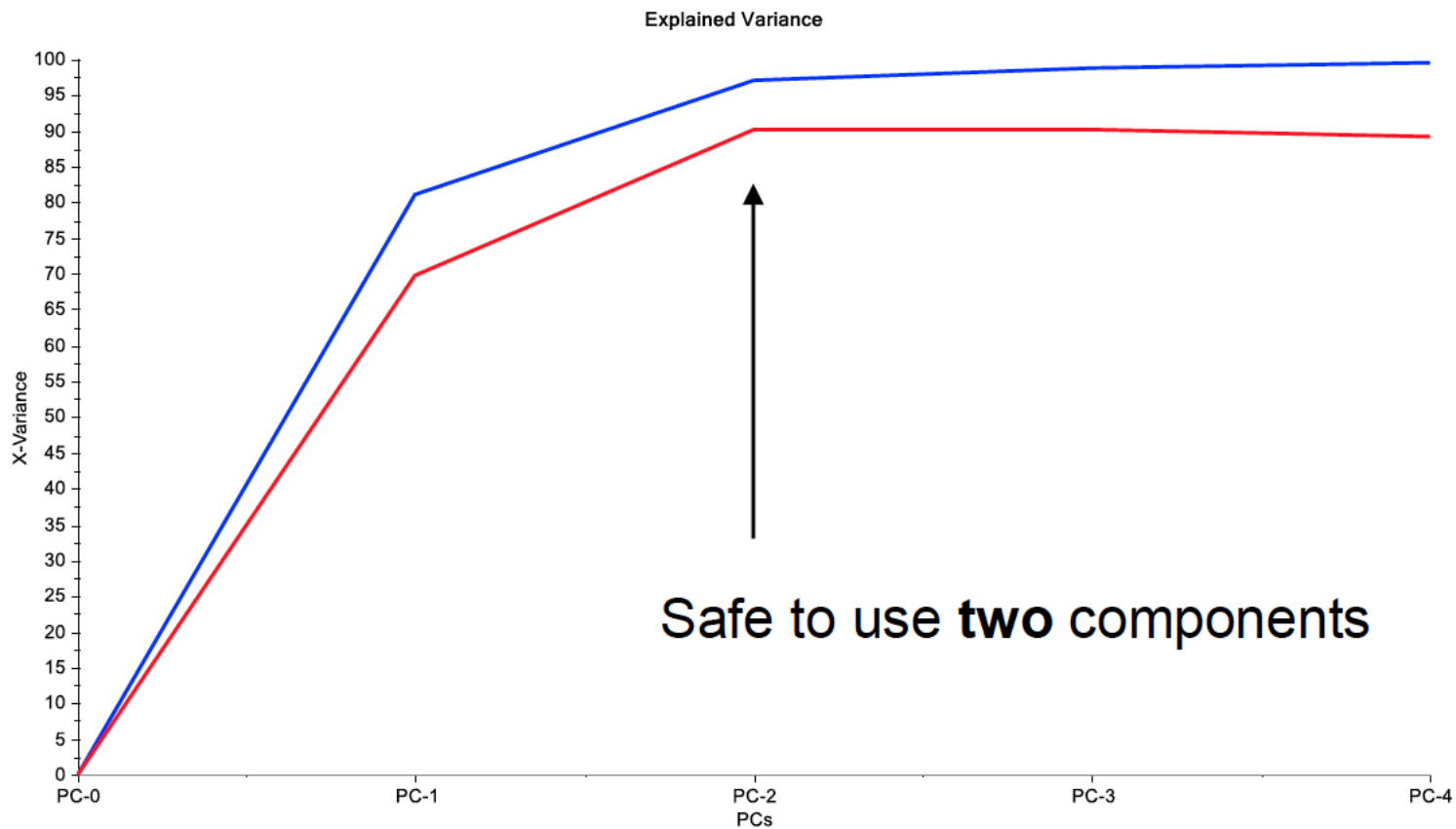


PCA - validation

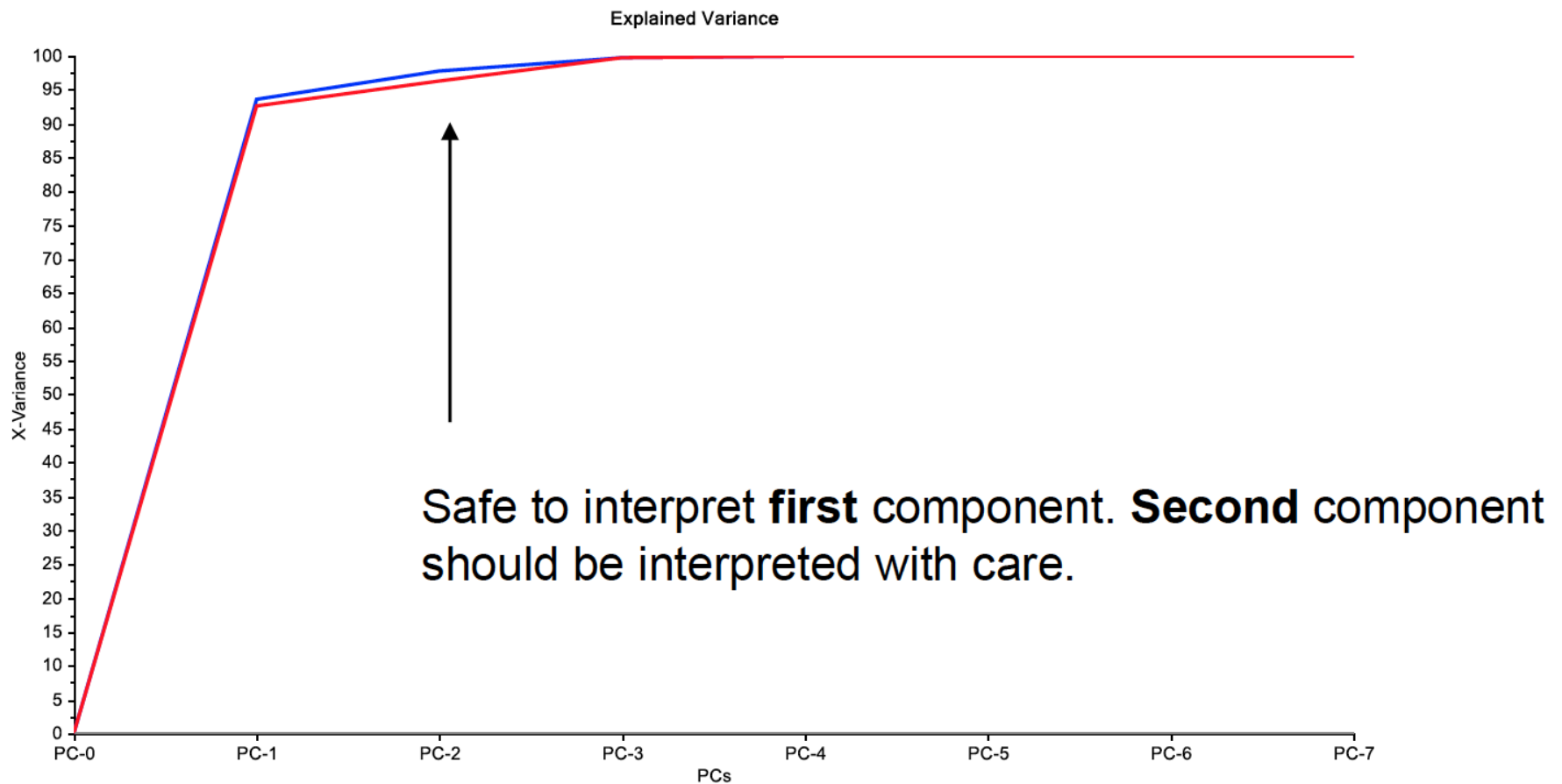
Validation is necessary to gain knowledge on **how many** components are **appropriate** for the model, i.e. how many components can be used for interpretation and further analysis.

- Use of *internal cross validation* in PCA
 - **K-fold** cross validation (number of folds / splits used)
 - **LOO** cross validation (“Leave-one-out”)
-

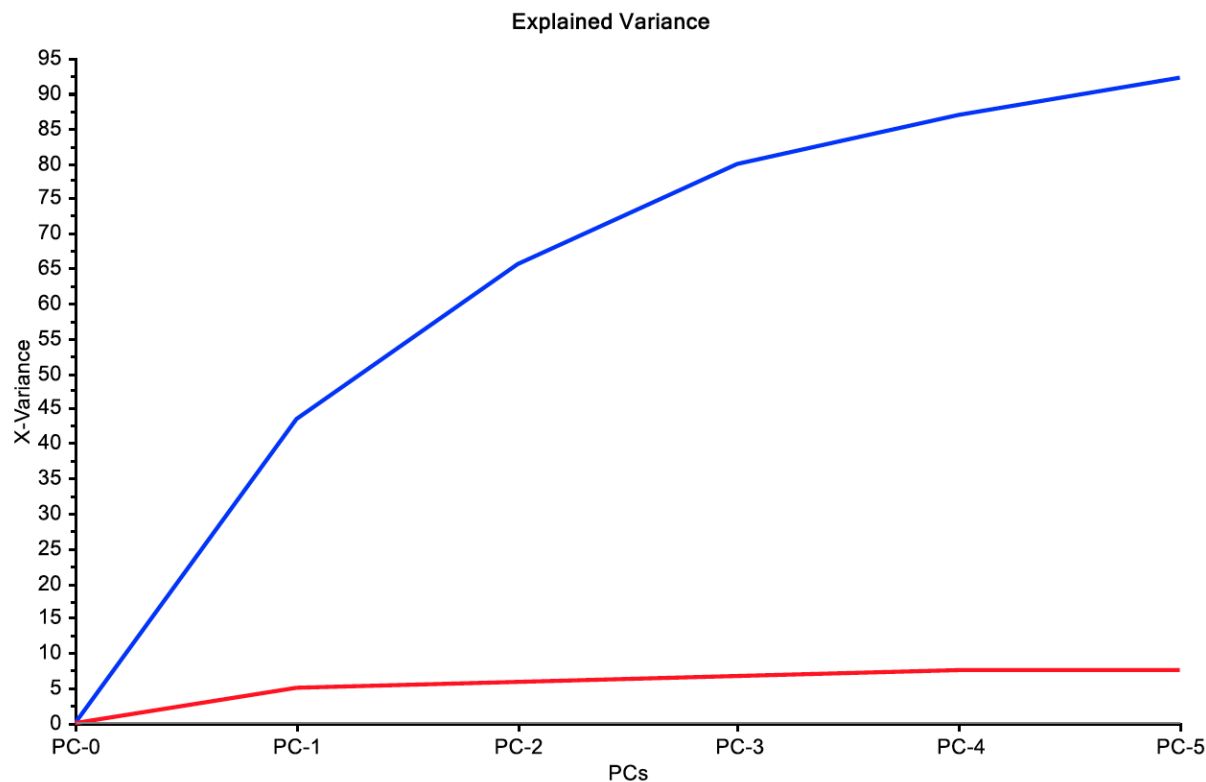
PCA - validation



PCA - validation



PCA - validation

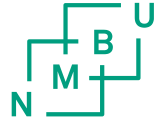


Poor model – may be a result of few objects that are very different from each other or overfitting

Hoggorm, HoggormPlot and examples

- **Hoggorm** package for multivariate statistics
 - GitHub: <https://github.com/olivertomic/hoggorm>
 - Read the Docs: <http://hoggorm.readthedocs.io/en/latest/>
- **HoggormPlot** package for convenient plotting of Hoggorm results
 - GitHub: <https://github.com/olivertomic/hoggormPlot>
 - Read the Docs: <http://hoggormplot.readthedocs.io/en/latest/>
- **Examples** of how to use Hoggorm illustrated in Jupyter notebooks
 - GitHub: <https://github.com/khliland/hoggormExamples>

Hoggorm



```
import hoggorm as ho
import hoggormplot as hopl
from sklearn import datasets

#import Iris data set

iris = datasets.load_iris()
X = iris.data
Y = iris.target

# Get the variabls
iris_varNames = list(iris.feature_names)

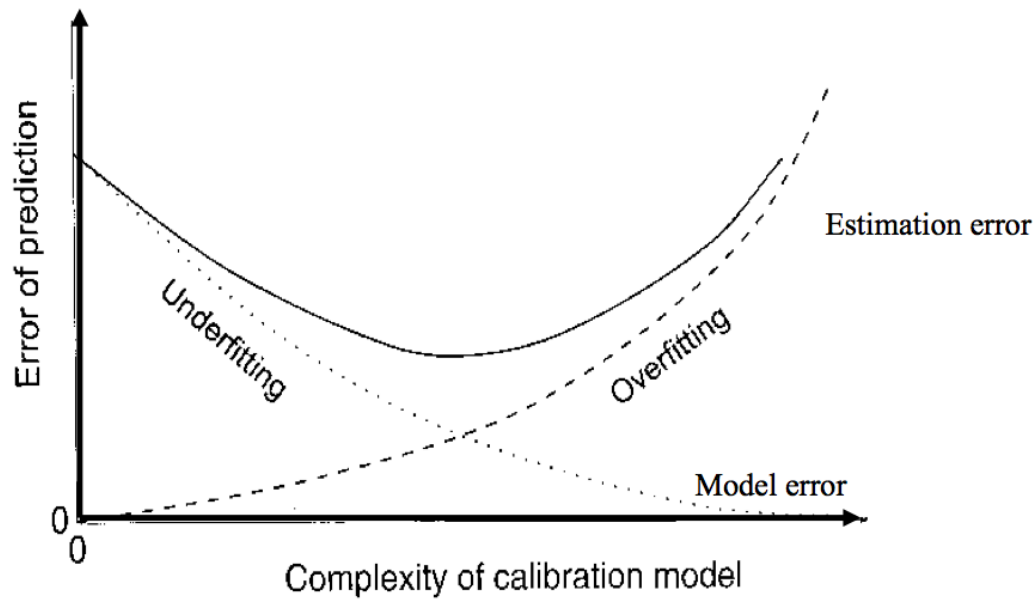
# Get the objects
iris_objNames = list(iris.target_names)

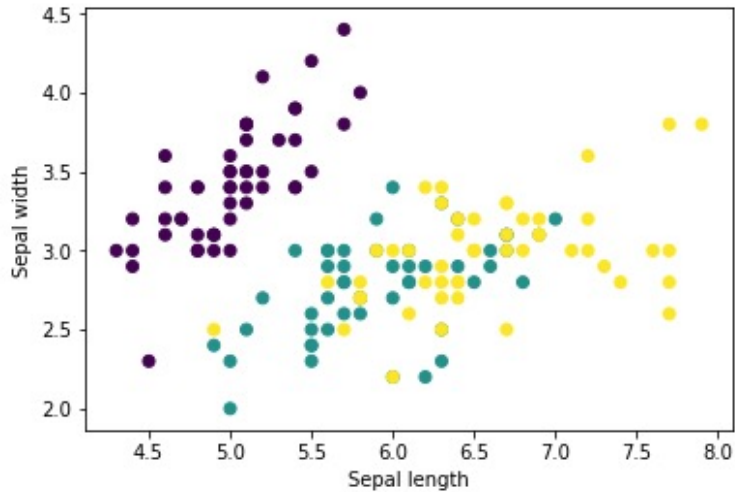
model_01 = ho.nipalsPCA(arrX=X, Xstand=True, cvType=["loo"], numComp=3)

hopl.plot(model_01, plots=[1, 2, 3, 6], line=True)
hopl.plot(model_01)
hopl.plot(model_01, plots=['scores', 'loadings', 'explainedVariance'])
```

HOW MANY COMPONENTS?

- «Use the number of components that do not produce overfitting of model»
- Explained variance for PCA
- Prediction error for PCR/PLS





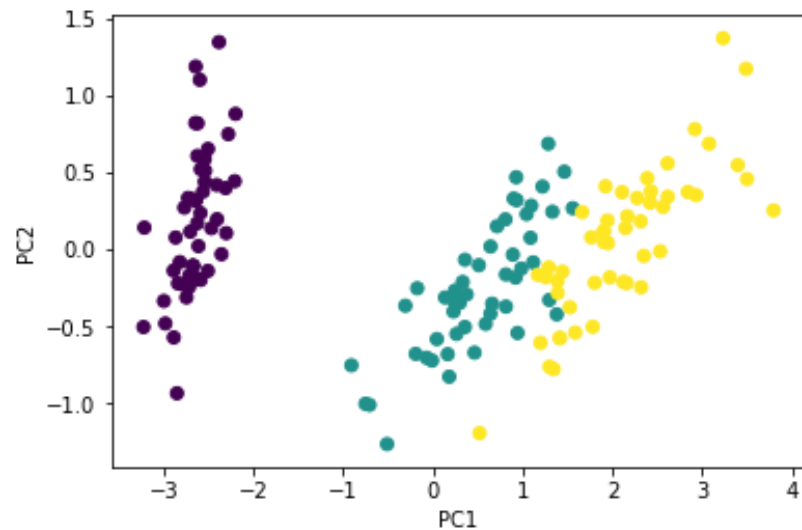
```
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.decomposition import PCA
```

```
# import the Iris data
iris = datasets.load_iris()
X = iris.data
y = iris.target
```

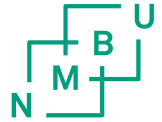
```
plt.figure()
plt.scatter(X[:,0],X[:,1],c=y)
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')
plt.show()
```

```
#simple PCA
```

```
X_reduced = PCA(n_components=3).fit_transform(iris.data)
plt.figure()
plt.scatter(X_reduced[:,1],X_reduced[:,2],c=y)
plt.xlabel('PC1')
plt.ylabel('PC2')
plt.show()
```

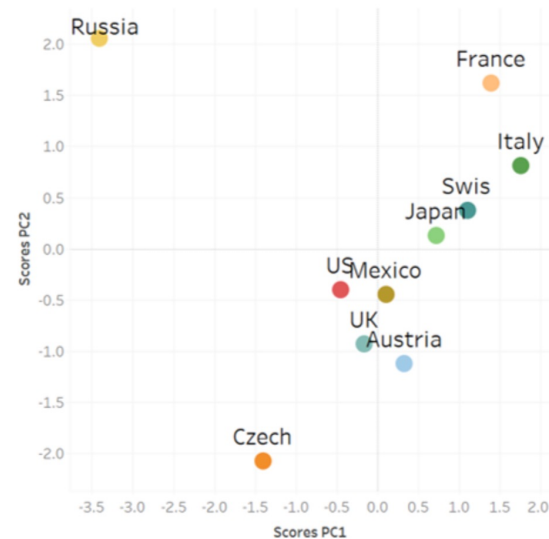
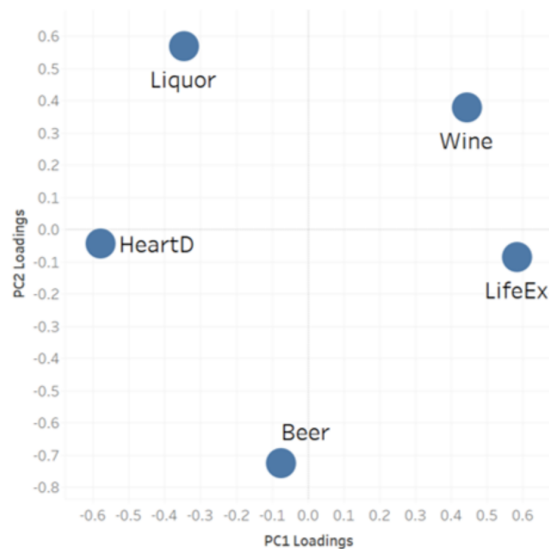


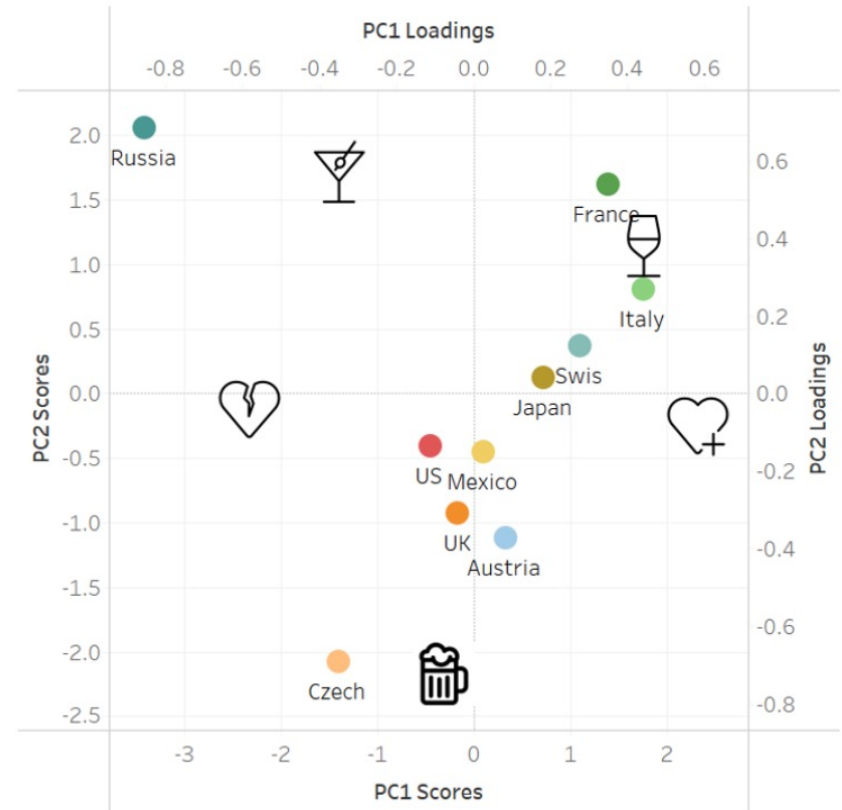
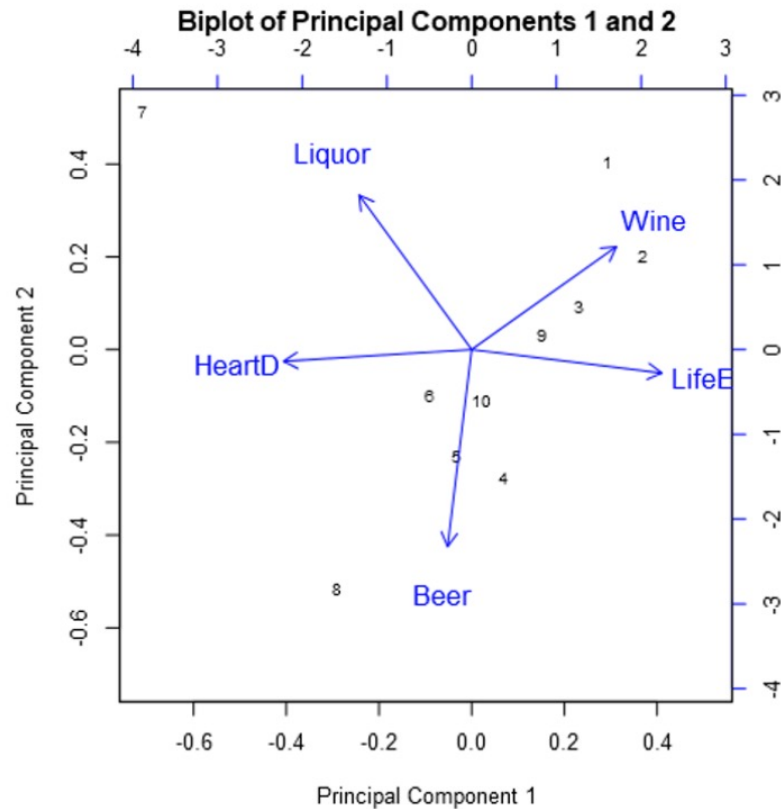
Causality ?



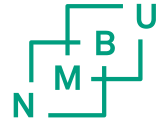
Data set from Time Magazines 1996

	Liquor	Wine	Beer	Life Ex	Heart D
France	2.5	63.5	40.1	78.0	61.1
Italy	0.9	58.0	25.1	78.0	94.1
Swis	1.7	46.0	65.0	78.0	106.4
Austria	1.2	15.7	102.1	78.0	173.0
UK	1.5	12.2	100.0	77.0	199.7
US	2.0	8.9	87.8	76.0	176.0
Russia	3.8	2.7	17.1	69.0	373.6
Czech	1.0	1.7	140.0	73.0	283.7
Japan	2.1	1.0	55.0	79.0	34.7
Mexico	0.8	0.2	50.4	73.0	36.4





- Does this mean that you live longer if you drink wine, shorter if you drink Liquor?
- The analysis is about correlation, not causality



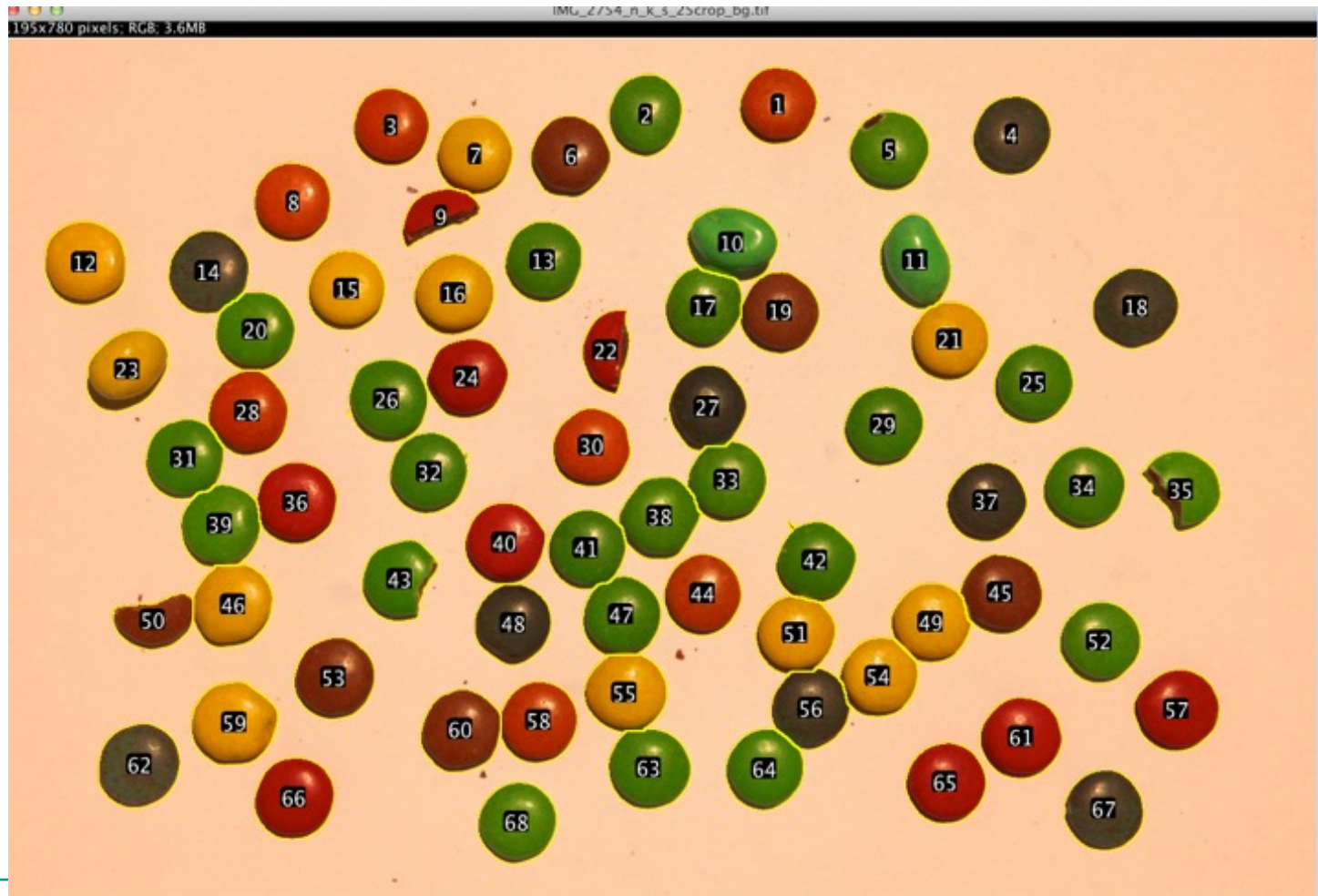
MULTIVARIATE ANALYSIS ON IMAGES

- Features in images are used in analysis
 - Matrix of Objects Shape, size, distribution etc
- Surface texture analysis
 - Image Surface texture



UNIVARIATE ANALYSIS ON IMAGES

- OBJECT analysis of chocolate

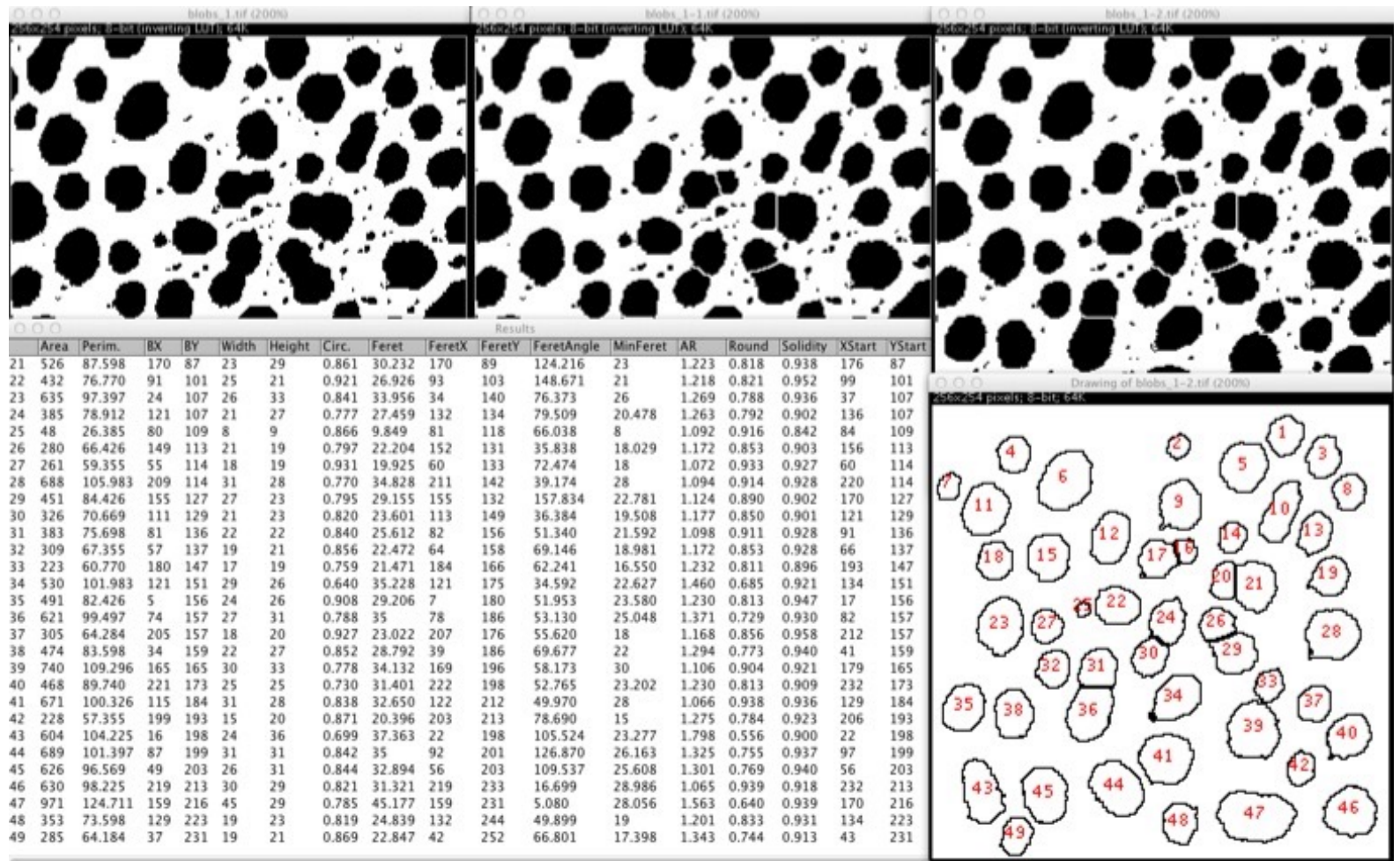


UNIVARIATE ANALYSIS ON IMAGES

- OBJECT analysis

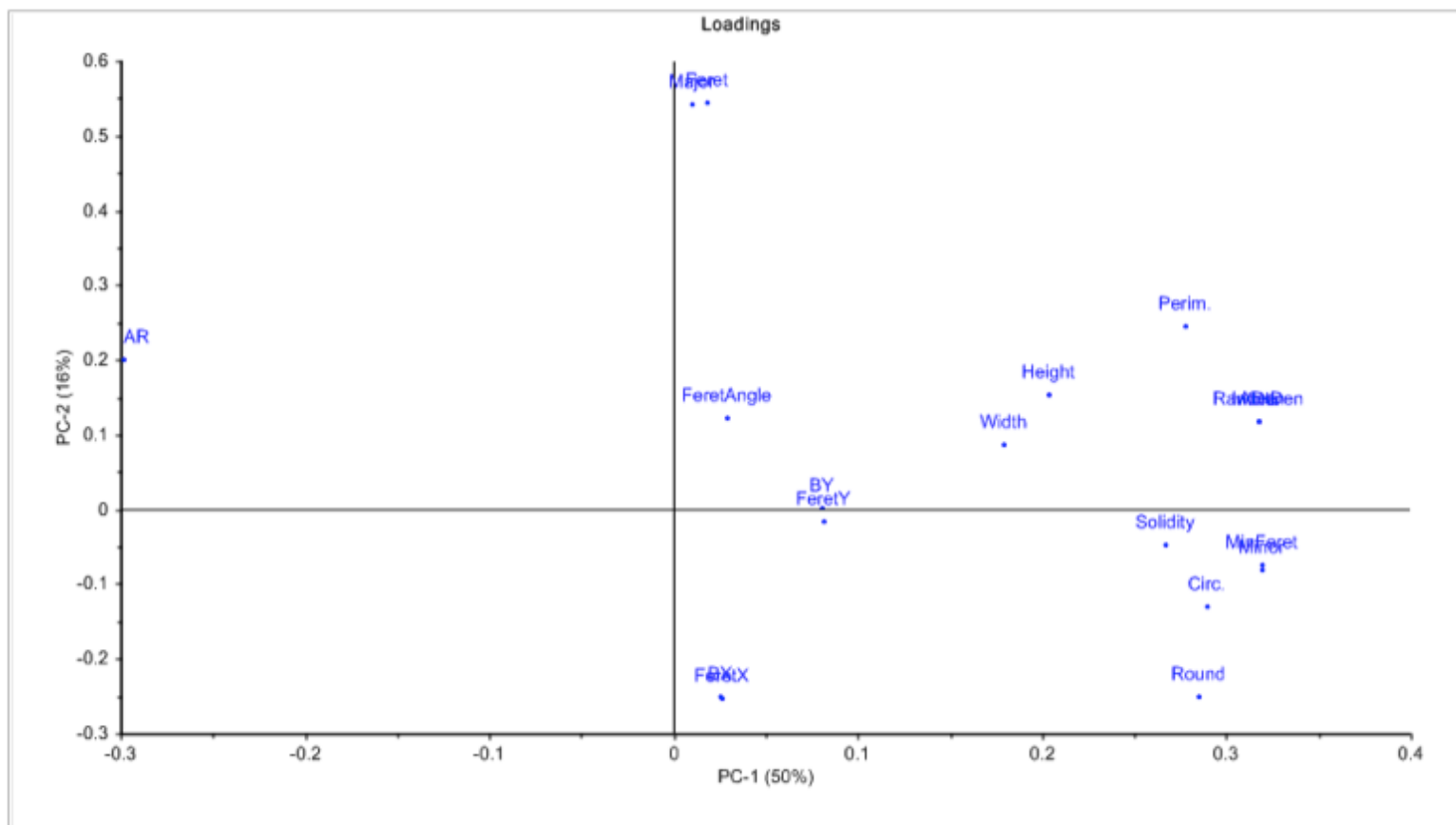


Figur 1. Sortert mhp «Roundness»



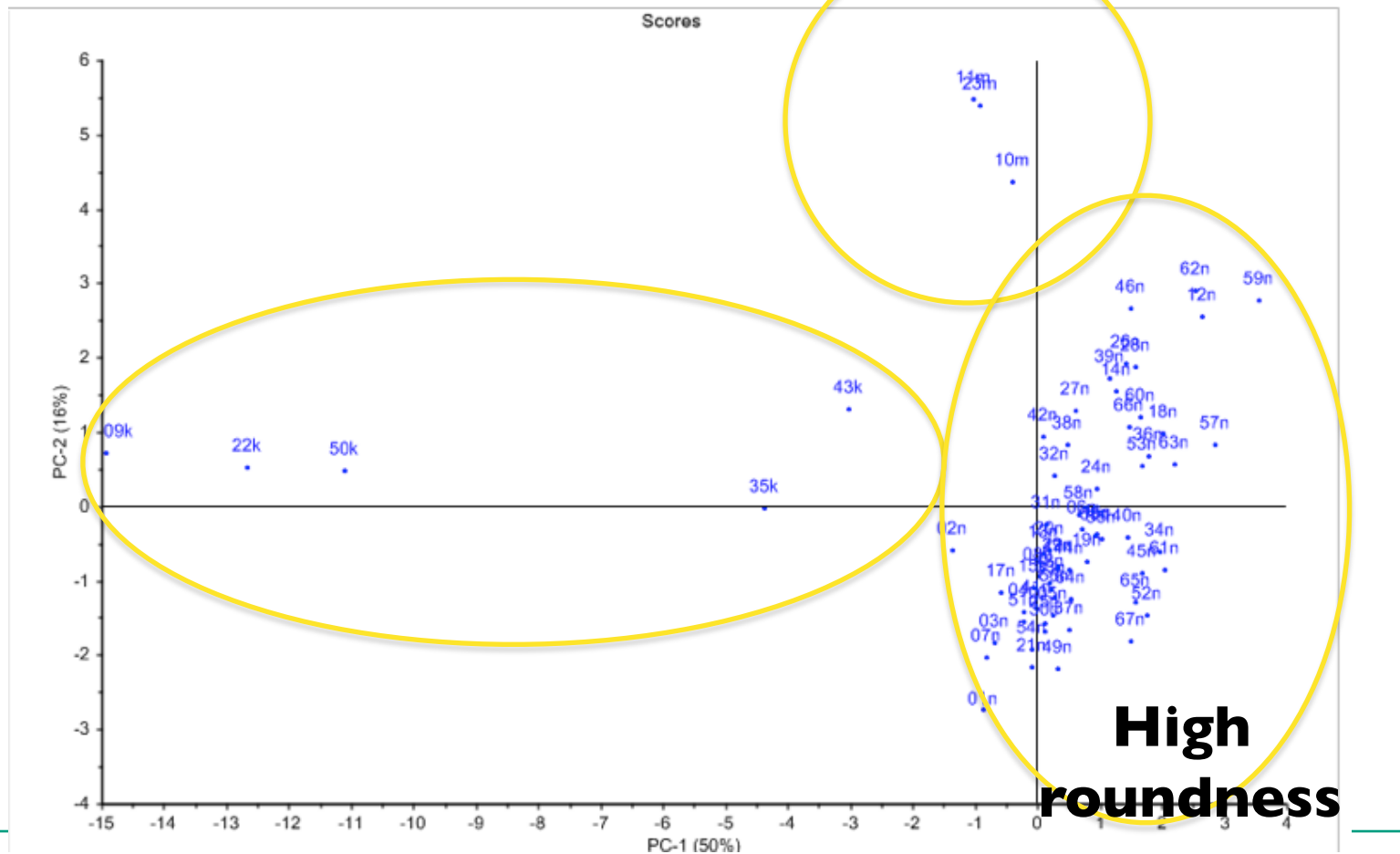
MULTIVARIATE ANALYSIS ON IMAGES

- PCA loadings on shape descriptors (area, roundness, ellipse etc)

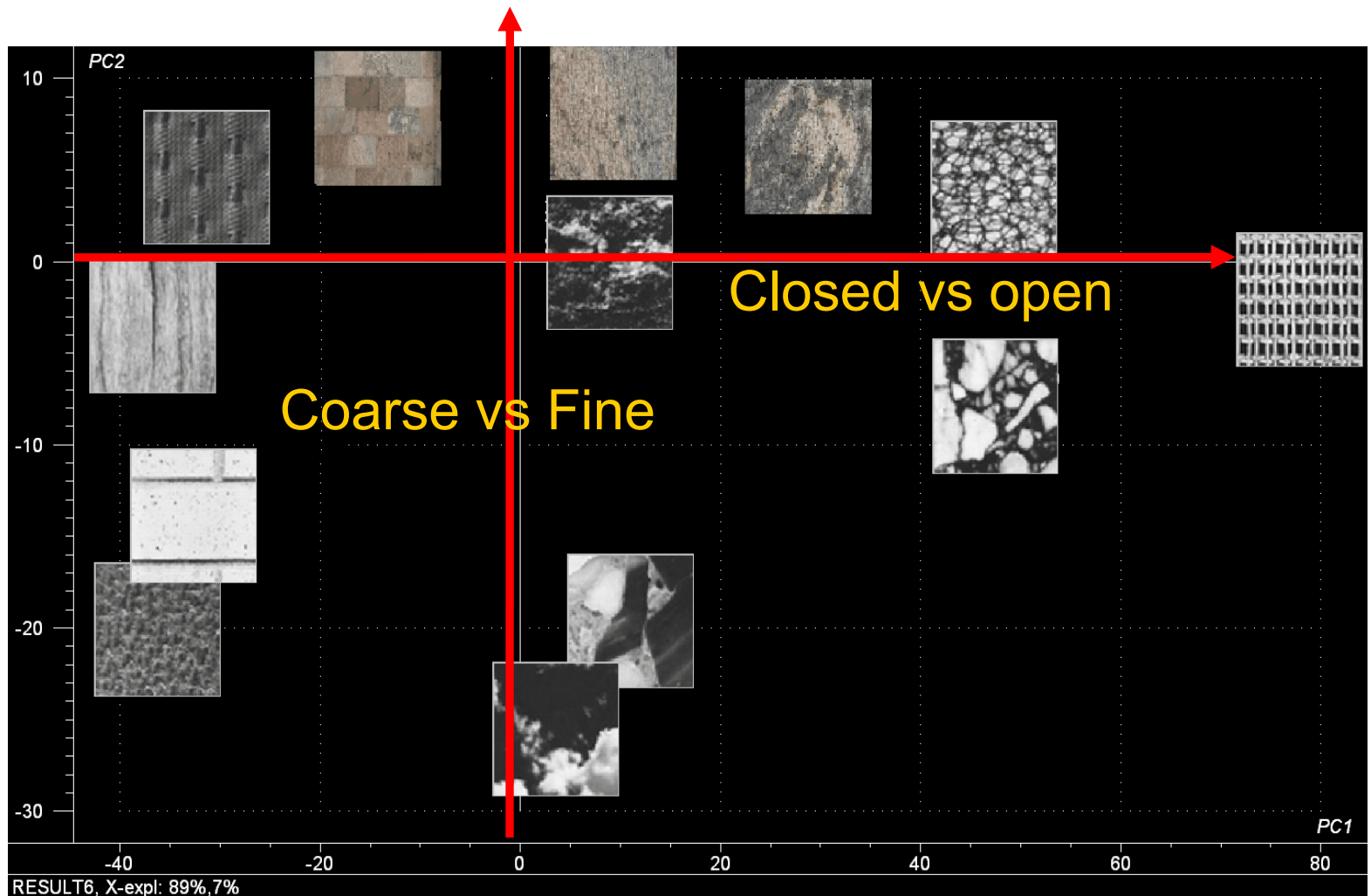
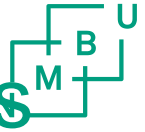


MULTIVARIATE ANALYSIS ON IMAGES

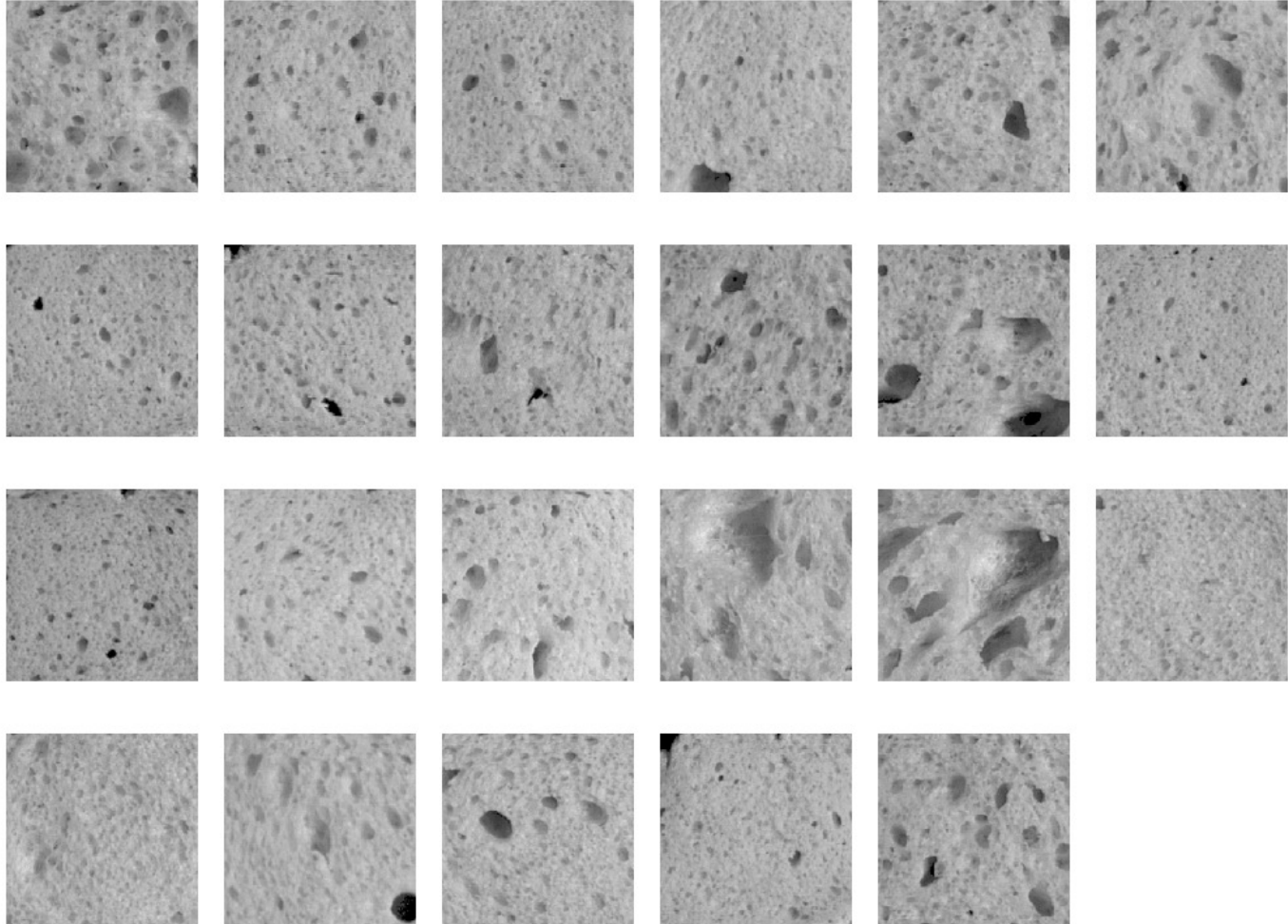
- PCA scores on shape descriptors (area, roundness etc)



Principal Components Analysis of textures



Baguette textures and modeling of sensory porosity



IMAGES MOUNTED AS POINTS IN PCA



Frost damage detection

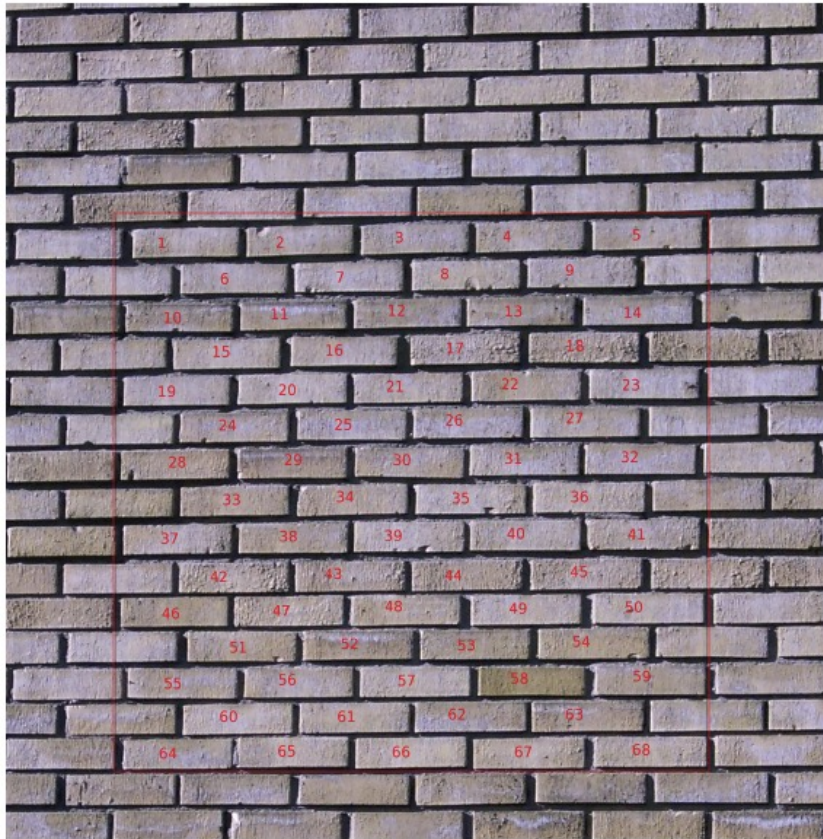


Fig 3: The original wall

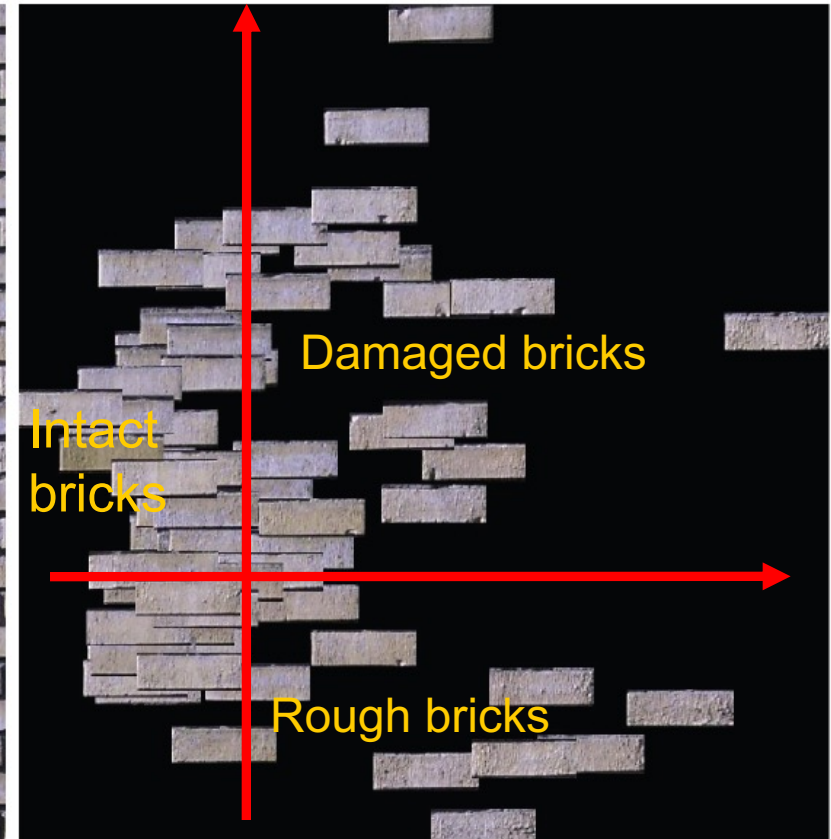
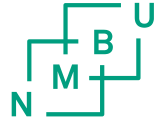
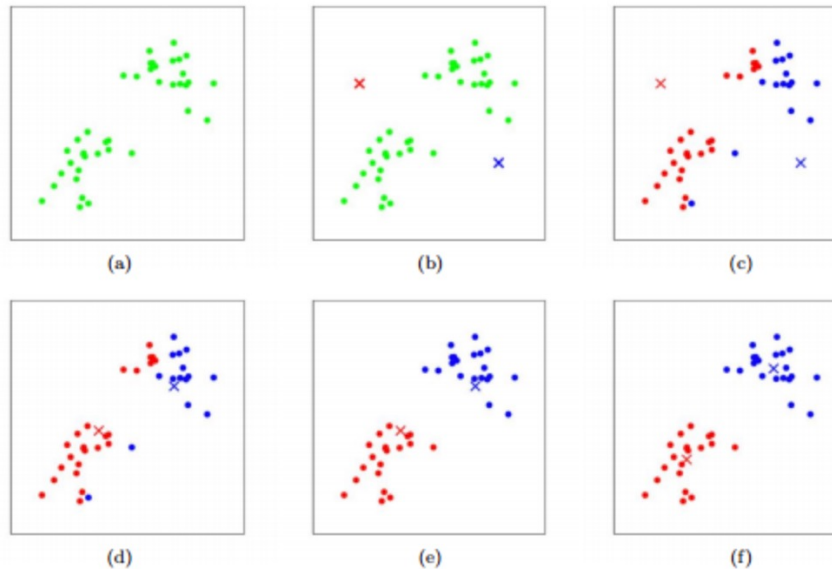


Fig 4: The bricks mounted in a PCA plot

Clustering



- K-means clustering
 - Separates samples into k groups of equal variance



– https://www.youtube.com/watch?v=4b5d3muPQmA&ab_channel=StatQuestwithJoshStarter

– <https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>