

Course 27411

Biological Data analysis and chemometrics

February 11th, 2008

"Missing in Action?" – please check-in

Anders C. Raffalt	s001831	Cecilie Kyrø	s043146
Julie K. Høgh	s020322	Louise M. Jørgensen	s050438
Lasse R. Bech	s021876	Fhayaz A. Sayyad	s060064
Mikkel T. Pleman	s022252	Emmanouil Papadakis	s060724
Christina H. Kærsgaard	s030513	Anders S. R. Ødum	s061888
Tina F. Rasmussen	s030526	Maria M. Bermejo	s071138
Ken Sejling	s030632	Carlos N. Bartolomé	s071139
Anders S. Laier	s031755	Asli Ozen	s071208
Rene J. Larsen	s031891	Silvia Bergamaschi	s073238
Jacob B. Reves	s032280	Michal Strejcek	s073676
Lars Poulsen	s032289	Gry H. Svendsen	s961360
Kasper Jensen	s042114	Casper T. Holst	s991097

Groups and Presentations

Group #1 (PCA):		Group #7 (Cluster Analysis):	
Lars Poulsen	s032289	Rene J. Larsen	s031891
Jacob B. Reves	s032280	Emmanouil Papadakis	s060724
Group #2 (MLR):		Group #8 (Classification I):	
Anders C. Raffalt	s001831	Anders S. R. Ødum	s061888
Casper T. Holst	s991097	Asli Ozen	s071208
Group #3 (PCR):		Group #9 (Classification II):	
Maria M. Bermejo	s071138	Silvia Bergamaschi	s073238
Carlos N. Bartolomé	s071139	Michal Strejcek	s073676
Group #4 (PLS-R):		Group #10 (PCO):	
Christina H. Kærsgaard	s030513	Fhayaz A. Sayyad	s060064
Tina F. Rasmussen	s030526	Kasper Jensen	s042114
Group #5 (Ridge regress.):		Anders S. Laier	s031755
Louise M. Jørgensen	s050438	Ken Sejling	s030632
Mikkel T. Pleman	s022252	Julie K. Høgh	s020322
Group #6 (Corresp. Analysis):		Lasse R. Bech	s021876
Cecilie Kyrø	s043146		
Gry H. Svendsen	s961360		

Lecture #2 Principal Component Analysis

Michael Adsetts Edberg Hansen
Assisting Professor, CMB/DTU
meh@bio.dtu.dk

February 11th, 2008

"Principal what?"

- Principal Component Analysis – PCA.
- "Invented" by Karl Pearson in 1901:

Pearson, K. (1901) On lines and planes of closest fit to systems of points in space. Philosophical Magazine (6) 2: 559-572.

- A.k.a. or closely related to:
 - Singular Value Decomposition (SVD)
 - Karhunen-Loève Expansion
 - Eigenvector Analysis
 - Latent Vector Analysis
 - Characteristic Vector Analysis
 - Hotelling Transformation

!!!! FIRST !!!!

- Lets recap from linear algebra

Matrix Inverse

- $A^{-1} A = A A^{-1} = I$

$$D D^{-1} = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 7 \end{bmatrix} \begin{bmatrix} \frac{1}{3} & 0 & 0 \\ 0 & \frac{1}{5} & 0 \\ 0 & 0 & \frac{1}{7} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = I$$

Properties

A^{-1} only exists if A is **square** ($n \times n$)

If A^{-1} exists then A is **non-singular** (invertible)

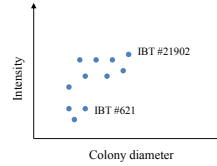
$(A B)^{-1} = B^{-1} A^{-1}$; $B^{-1} A^{-1} A B = B^{-1} B = I$

$(A^T)^{-1} = (A^{-1})^T$; $(A^{-1})^T A^T = (A A^{-1})^T = I$

Data as observations

- Samples as observations in a Multi/hyper-dimensional Space:

- Objects are a collection of features.
- Features are dimensions.
- Objects are points in a multidimensional space.



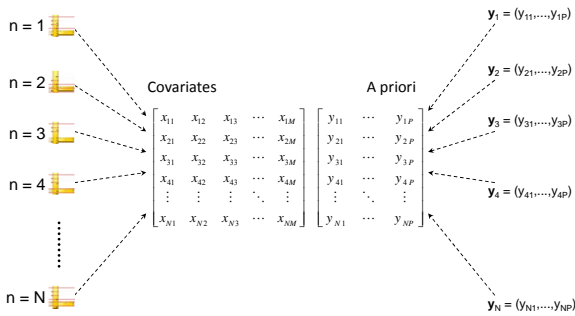
- Mathematical notation

- N is the number of observations
- M is the number of variables/features

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1M} \\ x_{21} & x_{22} & \dots & x_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NM} \end{bmatrix}$$

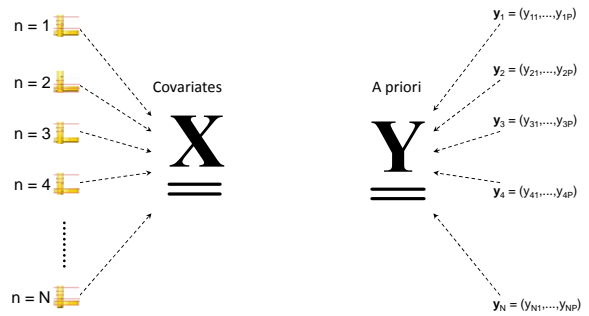
Data as observations

- For each sample



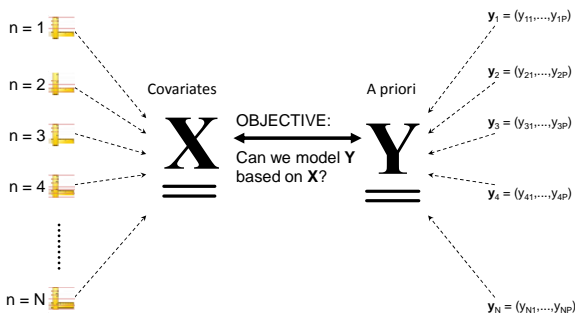
Data as observations

- For each sample



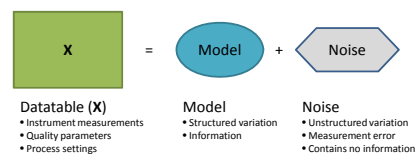
Data as observations

- For each sample



Principal Component Analysis (PCA)

- Projection method
- Exploratory data analysis
- Extract information and remove noise
- Reduce dimensionality / Compression
- Classification and clustering

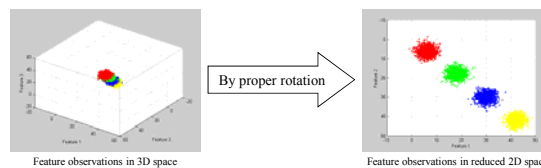


PCA in a nutshell

- First...
... The illustrational/intuitive approach:

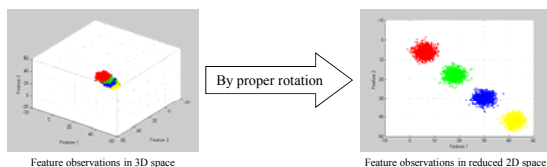
Projection of data

- Linear transformation



Projection of data

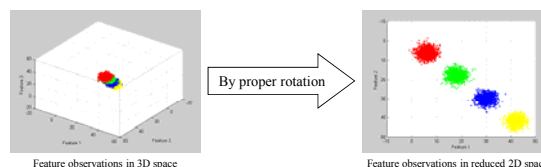
- By proper linear transformation



- The PCA approach:
 - Rotation according to maximum variance in data.

Projection of data

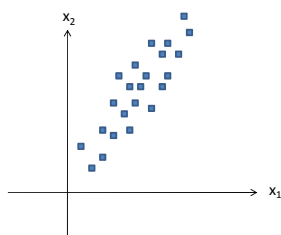
- By proper linear transformation



- The PCA approach:
 - Rotation according to maximum variance in data.
- Fisher approach (**later lecture**):
 - Rotation according to maximum discrimination between groups.

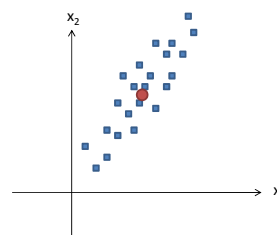
Principal Component Analysis

- For a given dataset:



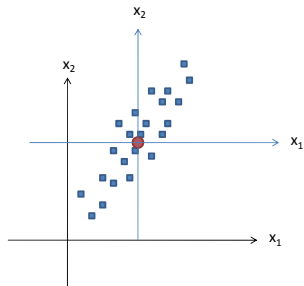
Principal Component Analysis

- Calculate the centroid (= "mean in all directions"):



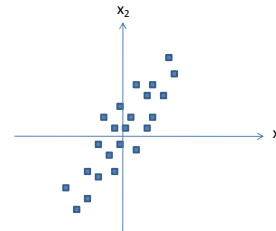
Principal Component Analysis

- Shift the grid to the centroid:



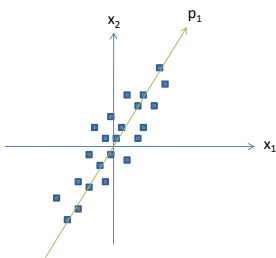
Principal Component Analysis

- Take this as our new coordinate system:



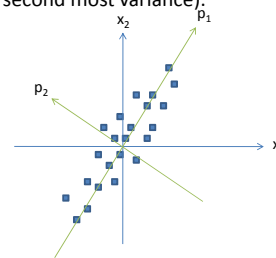
Principal Component Analysis

- Calculate the direction in which the variance is maximal:



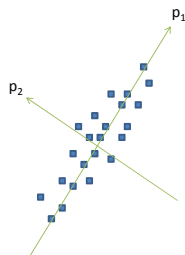
Principal Component Analysis

- And repeat this for each next perpendicular axis (direction with second most variance):



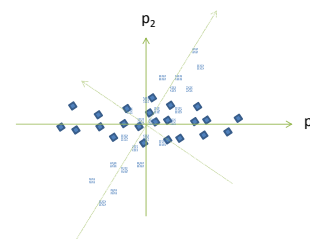
Principal Component Analysis

- Leaving us with a rotated grid:



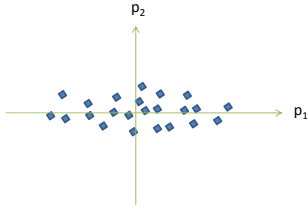
Principal Component Analysis

- Which we can rotate to a "normal" position:



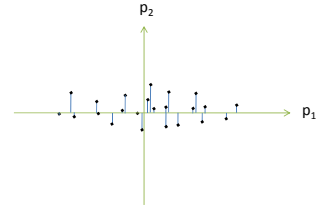
Principal Component Analysis

- Showing us maximal variance:



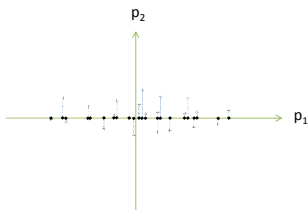
Principal Component Analysis

- We can also use this to reduce the complexity of the data set:



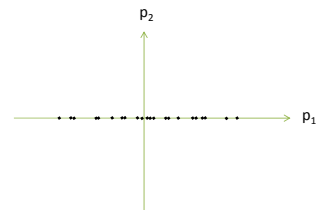
Principal Component Analysis

- By eliminating a number of axis by projection of the points:



Principal Component Analysis

- In this example moving from two...:

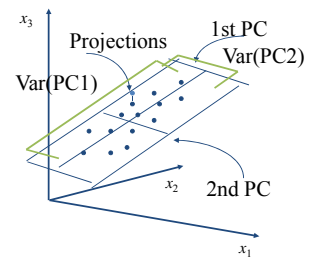


Principal Component Analysis

- ... to one dimensional data points:

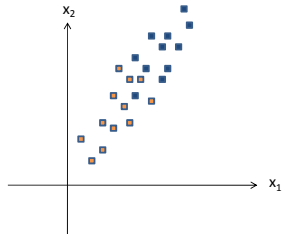


In 3D



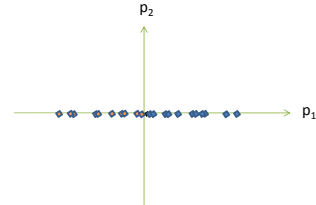
Principal Component Analysis

- Assuming variation equals species diversity...



Principal Component Analysis

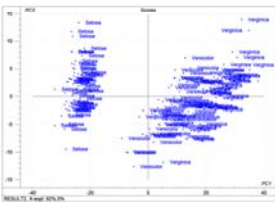
- ... the first PCA depicts this information...



Scores and Loadings

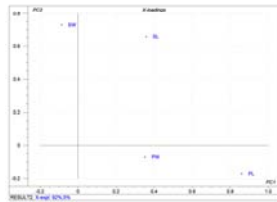
Scores

- Map of samples
- Displays distribution of samples in the new space defined by the PC's



Loadings

- Map of variables
- Shows how the original variables are related to the PC's



Principal Component Analysis

- We center the data.
 - But what about scaling?
- Especially when the covariates (columns) represent different scales (i.e. comparing apples and bananas) its important to divide by the standard deviation!

PCA in a nutshell

- Next...
 - ... lets describe the illustrations in mathematical terms:

Some (boring) definitions – 1 of 2

- PCA is mathematically defined as an **orthogonal linear transformation** that transforms the data to a new coordinate system **such that the greatest variance by any projection of the data comes to lie on the first coordinate** (called the first principal component), the second greatest variance on the second coordinate, and so on.

Some (boring) definitions – 2 of 2

- PCA can be used for **dimensionality reduction** in a data set by **retaining those characteristics of the data set that contribute most to its variance**, by keeping lower-order principal components and ignoring higher-order ones.
- Such low-order components often contain the “most important” aspects of the data.
 - i.e. **assuming that information is equal variation!**
- However, depending on the data this may (obviously) not always be the case.

PCA in a nutshell

- Finally... Lets wrap it up in math:

Principal Component Analysis

- The *i*'th principal component of \underline{X} is the projection

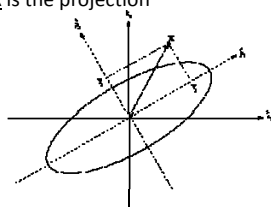
$$\underline{Y}_i = \underline{p}_i^T \underline{X}$$

- The vector \underline{Y}

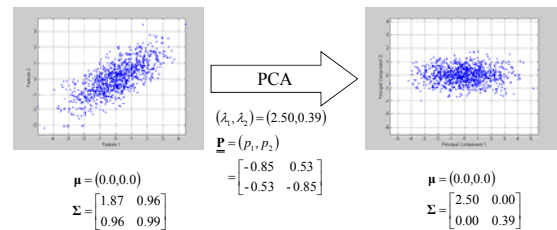
$$\underline{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_k \end{pmatrix} = \underline{P}^T \underline{X}$$

is called the vector of principal components.

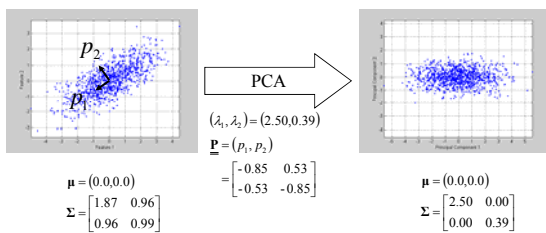
$$\text{cov}(\underline{Y}) = \text{cov}(\underline{P}^T \underline{X}) = \underline{P}^T \underline{\Sigma} \underline{P} = \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_k \end{pmatrix}$$



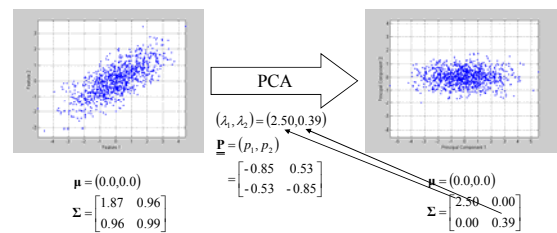
Principal Component Analysis



Principal Component Analysis



Principal Component Analysis



Principal Component Analysis

- From the PCA we may extract the set of linear combinations that explains the most variation

$$\frac{\lambda_1 + \dots + \lambda_m}{\lambda_1 + \dots + \lambda_m + \dots + \lambda_k} \geq q$$

- And hereby condense and reduce the dimensionality of the featurespace.
- From the example before we see, that the linear combinations explain

$$\underline{\mathbf{P}} = (p_1, p_2) = (86.3\%, 13.7\%)$$

of the total variance.

Output

- Loadings
 - The weights
- Scores
- Plots
 - Loadings plot
 - Scores plot
 - Biplot

Pros and Cons

- Positives**
 - Can deal with large data sets.
 - There weren't done any assumptions on the data. This method is general and may be applied to any data set.
- Negatives**
 - Nonlinear structure is invisible to PCA
 - The meaning of features is lost when linear combinations are formed
- Still!!!!**
 - Nonlinear PCA's exist (so-called kernel methods)
 - Sparseness or supervised projections can be introduced to emphasize important features

Course 27411 – Exercises 11/2 2008

PCA Exercise 1, Fisher iris data

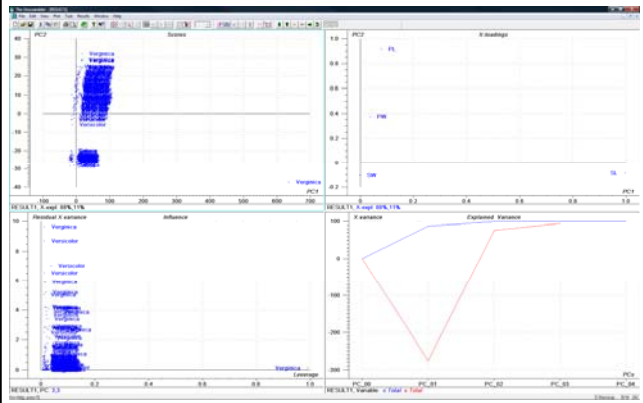
The Fisher iris data set is classical¹. There are 150 objects, 50 iris setosa, 50 iris versicolor and 50 iris virginica. The flowers of these 150 plants have been measured by a ruler. The variables are sepal length, sepal width, petal length and petal width (4 variables).

The original hypothesis was that *I. versicolor* was a hybrid of the two other species i.e. *I. setosa* x *I. virginica*. *I. setosa* is a diploid species; *I. virginica* is a tetraploid and *I. versicolor* is a hexaploid.

- Start the Unscrambler and open the help function. Lookup PCA and read about the different options and plots provided by the software.
In general: Use the help function before asking! ☺ To get help on the different plots mark the plot and press F1.
- Import the excel file fishiris.xls to UNSCRAMBLER and save the UNSCRAMBLER file
- Usually you should examine the raw data first, but in this exercise try to perform a PCA, with "leverage correction" and with centering. Look at the four standard plots.
 - How many principal components would you need and what does the PC1 describe?
 - How many % of the variations is described by the first two PCs?
 - Do you see anything "wrong" in the score-plot?
 - If so what is wrong and what should be done?
 - Try to correct it!
- Does it help to standardize the data matrix and perform a new PCA (try it!).
- How many PCs are needed to explain 70, 75% or 90% of the variance?
- What does PC1, PC2 and PC3 show?
- Are any variables more discriminative than others? Are any variables dispensable?
- Can you see the presupposed classes? Any class overlap?

¹ Fisher, R.A. (1936). "The Use of Multiple Measurements in Taxonomic Problems". *Annals of Eugenics* 7: 179-188.
² Edgar Anderson (1935). "The Vase of the Gauge Peninsula". *Bulletin of the American Iris Society* 5: 2-5.

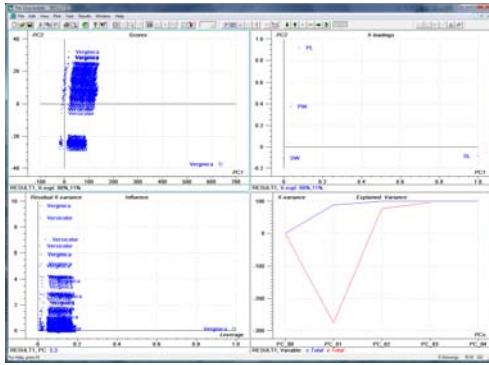
Exercise



Exercise



Exercise



Exercise

