



Principal Component Analysis

Jingwen Li

12.08.2024



Contents

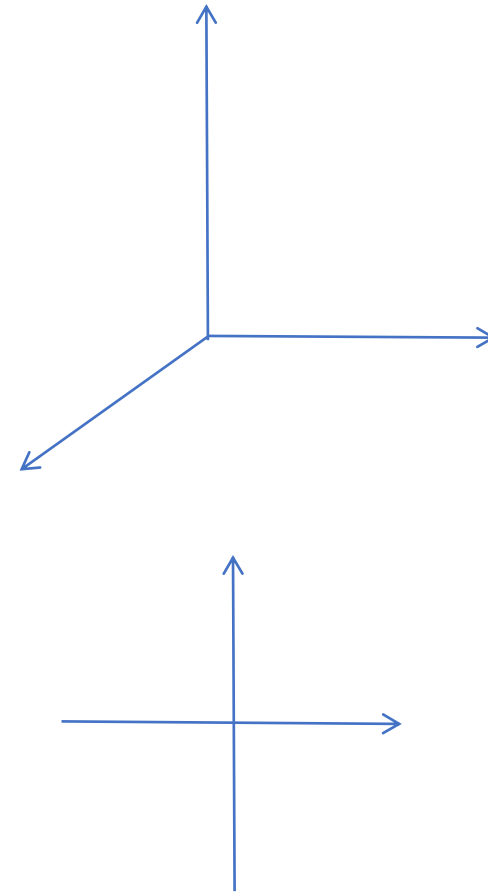
- 1 Definition
- 2 History
- 3 Process
- 4 Implications in data science
- 5 Alternative technique

Definition of PCA

Part.01

Definition of PCA

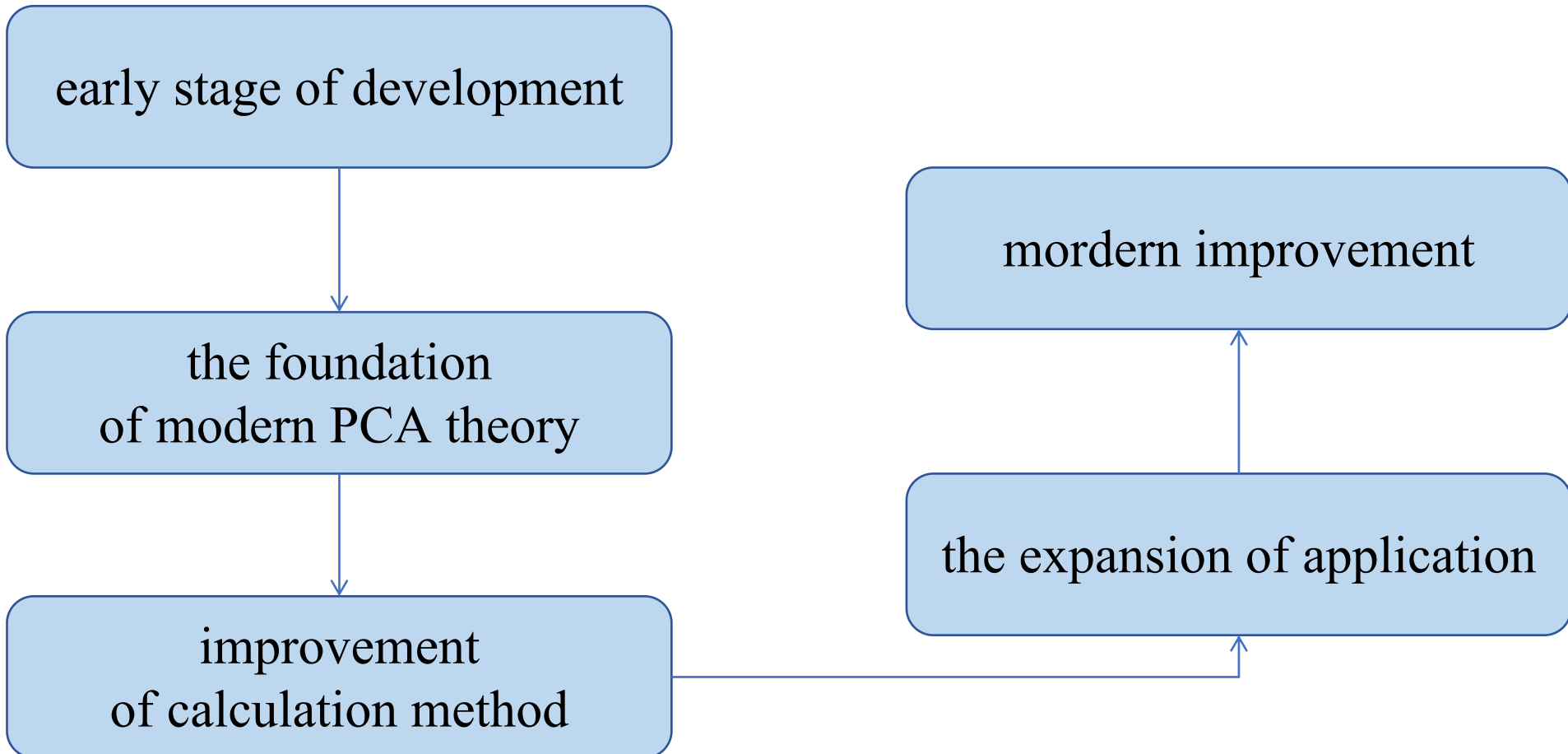
- Principal Component Analysis(PCA) is a statistical technique used for dimensionality reduction, which transforms a high-dimension data set into a low-dimension one while retaining as much of the original information as possible.



History of PCA

Part.02

History



Process of PCA

Part.03

Process of PCA

- Data standardization

$$\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1m} \\ x_{21} & x_{22} & \cdots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nm} \end{bmatrix}$$

$$\textcircled{1} \mu_j = \frac{1}{n} \sum_{i=1}^n$$

$$\textcircled{2} \sigma_j = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_{ij} - \mu_j)^2}$$

$$\textcircled{3} X_{\text{standardized},ij} = \frac{X_{ij} - \mu_j}{\sigma_j}$$

```
data matrix after data standardized(X_standardized):
  sepal length (cm)  sepal width (cm)  petal length (cm)  petal width (cm)
0      -0.900681      1.019004      -1.340227      -1.315444
1      -1.143017      -0.131979      -1.340227      -1.315444
2      -1.385353      0.328414      -1.397064      -1.315444
3      -1.506521      0.098217      -1.283389      -1.315444
4      -1.021849      1.249201      -1.340227      -1.315444
...      ...      ...      ...      ...
145      1.038005      -0.131979      0.819596      1.448832
146      0.553333      -1.282963      0.705921      0.922303
147      0.795669      -0.131979      0.819596      1.053935
148      0.432165      0.788808      0.933271      1.448832
149      0.068662      -0.131979      0.762758      0.790671

[150 rows x 4 columns]
```


Process of PCA

- Calculate covariance matrix

$$\Sigma = \frac{1}{n-1} X_{\text{Standardized}}^T X_{\text{Standardized}}$$

Process of PCA

- Calculate covariance matrix

$$X_{\text{standardized}} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1m} \\ x_{21} & x_{22} & \cdots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nm} \end{bmatrix} \quad X_{\text{Standardized}}^T = \begin{bmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1m} & x_{2m} & \cdots & x_{nm} \end{bmatrix}$$

$$C = X_{\text{Standardized}}^T \times X_{\text{Standardized}}$$

$$c_{ij} = \sum_{k=1}^n x_{ki} \times x_{kj}$$

Process of PCA

- Calculate eigenvalues and eigenvectors

$$\Sigma v = \lambda v$$

$$\det(\Sigma - \lambda I) = 0$$

$$(\Sigma - \lambda_i I)v_i = 0$$

```
eigenvalues(Eigenvalues):
[2.93808505 0.9201649 0.14774182 0.02085386]

eigenvectors(Eigenvectors):
```

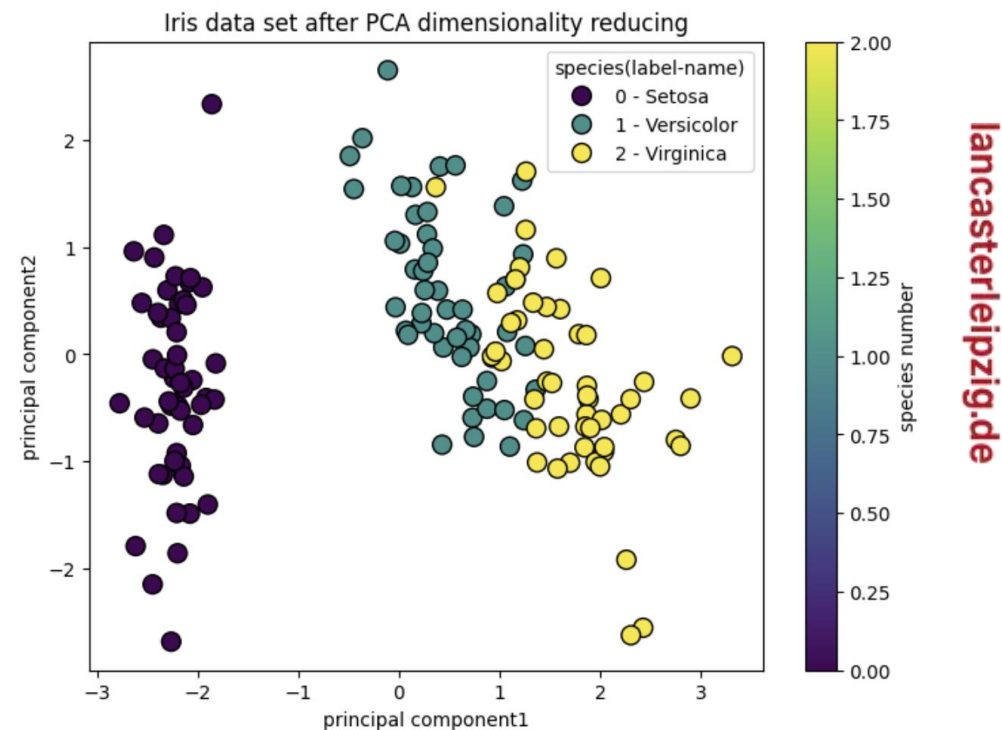
	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
sepal length (cm)	0.521066	-0.377418	-0.719566	0.261286
sepal width (cm)	-0.269347	-0.923296	0.244382	-0.123510
petal length (cm)	0.580413	-0.024492	0.142126	-0.801449
petal width (cm)	0.564857	-0.066942	0.634273	0.523597

Process of PCA

- Select principal component
- Construct transformation matrix
- Data transformation and dimensionality reduction

$$W = \begin{bmatrix} v_1 & v_2 & \cdots & v_k \end{bmatrix}$$

$$X' = W^T X_{\text{standardized}}$$



Implications in data science

Part.04

Implications in data science

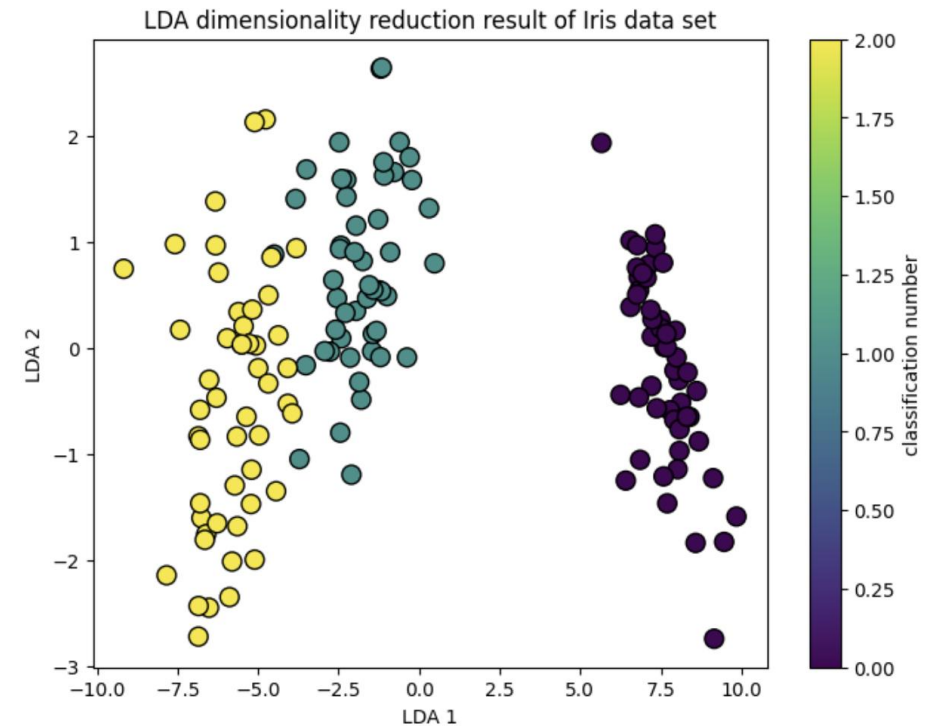
- Dimensionality Reduction
- Noise reduction
- Data visualization
- Feature extraction

Alternative technique

Part.05

LDA(Linear Discriminant Analysis)

- Comparison with PCA, LDA is supervised and uses category labels to find directions that best differentiate between different categories. This means that LDA generally performs better than PCA in classification tasks.



Thank you for your time and patience.