Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○

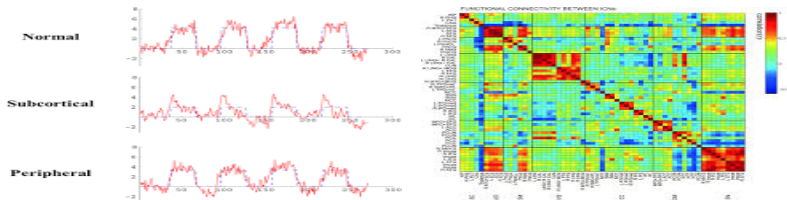t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

Dimension reduction with PCA and t-SNE

Motivation
●○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# High dimensional data in data analysis?



Words embeddings in NLP

Motivation
○●○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

## High dimensional data in data analysis?



Brain activity

Motivation
○○●○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

## High dimensional data in data analysis?

Challenges ?

Motivation
00000000

Principal component analysis
0000000000000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

## High dimensional data in data analysis?

- Challenges ?
  - Visualize
  - Group in relevant clusters
- Difficult with high dimensional data!
- A classical dimension reduction approach Principal Component Analysis

# Dimension reduction?



Dimension reduction without loss of information?

Motivation
○○○○○●○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

## Dimension reduction



Dimension reduction without loss of information?

**Motivation**
○○○○○○●○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

## Dimension reduction

### How can we reduce dimension to separate observations?

Two examples in this lecture

- A linear approach : Principal Component Analysis (PCA)
- A non linear alternative : t-distributed stochastic neighbourhood embedding (t-SNE)

Motivation
○○○○○○○●

Principal component analysis
○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

## Dimension reduction

### PCA vs t-SNE?

- Principal Component analysis (PCA)
  - preserves the global structure of data.
  - maps all the clusters as a whole
  - potential applications : noise filtering, feature extractions, stock market predictions, and gene data analysis.
- t-distributed stochastic neighbourhood embedding (t-SNE)
  - preserves the local structure of data
  - potential applications : music analysis, bioinformatics, and biomedical signal processing

Motivation
○○○○○○○○

Principal component analysis
●○○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis

Principle : find a linear projection on a low-dimensional space



How can we find the low-dimensional space $H$?

Motivation
00000000

Principal component analysis
0●000000000000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

# Principal Component Analysis
Practical examples with Python

Let us consider a data set describing tree kinds of leafs coming from the website : https://archive.ics.uci.edu/ml/datasets/Leaf

### A toy example

- Leaf data set : describe tree kinds of leafs coming from the website : https://archive.ics.uci.edu/ml/datasets/Leaf
  Focus on the two following variables
  - Elongation : maximal normalized distance between a point of the leaf and its boundary
  - Isoperimetric factor : ratio between the area and the square of the perimeter of the leaf

Motivation
○○○○○○○○

Principal component analysis
○○●○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis

Practical examples with Python

### Leaf Data Set
*Download*: Data Folder, Data Set Description

**Abstract**: This dataset consists in a collection of shape and texture features extracted from digital images of leaf specimens originating from a total of 40 different plant species.

| Data Set Characteristics: | Multivariate | Number of Instances: | 340 | Area: | Computer |
|---|---|---|---|---|---|
| Attribute Characteristics: | Real | Number of Attributes: | 16 | Date Donated | 2014-02-24 |
| Associated Tasks: | Classification | Missing Values? | N/A | Number of Web Hits: | 88647 |

**Source:**

This dataset was created by Pedro F. B. Silva and André R. S. Marçal using leaf specimens collected by Rubim Almeida da Silva at the Faculty of Science, University of Porto, Portugal.

**Data Set Information:**

For further details on this dataset and/or its attributes, please read the 'ReadMe.pdf' file included and/or consult the Master's Thesis 'Development of a System for Automatic Plant Species Recognition' available at [Web Link].

**Attribute Information:**

1. Class (Species)
2. Specimen Number
3. Eccentricity
4. Aspect Ratio
5. Elongation
6. Solidity
7. Stochastic Convexity
8. Isoperimetric Factor
9. Maximal Indentation Depth
10. Lobedness
11. Average Intensity
12. Average Contrast
13. Smoothness
14. Third moment
15. Uniformity
16. Entropy

Motivation
00000000

Principal component analysis
000●000000000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

## Principal Component Analysis
Practical examples with Python

```
import numpy as np
leaf =
np.loadtxt('/home/marianne/Downloads/leaf.csv',
delimiter=',')
import pandas as pd
df_leaf=pd.DataFrame(np.array([leaf[:,4],leaf[:,7]]).T,colu
'Isoperimetric factor'])
df_leaf
```

Motivation
○○○○○○○○

Principal component analysis
○○○○●○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis
Practical examples with Python

|      | Elongation | Isoperimetric factor |
| ---- | ---------- | -------------------- |
| 0    | 0.32396    | 0.835920             |
| 1    | 0.36116    | 0.798670             |
| 2    | 0.38998    | 0.808120             |
| 3    | 0.35376    | 0.816970             |
| 4    | 0.44462    | 0.754930             |
| ...  | ...        | ...                  |
| 335  | 0.81725    | 0.125230             |
| 336  | 0.75319    | 0.136860             |
| 337  | 0.78147    | 0.135030             |
| 338  | 0.71532    | 0.157470             |
| 339  | 0.85409    | 0.078376             |

Motivation
00000000

Principal component analysis
00000●00000000000000

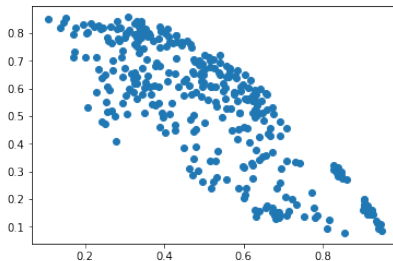t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

# Principal Component Analysis
Practical examples with Python

```
import matplotlib.pyplot as plt
fig, ax = plt.subplots()
ax.scatter(leaf[:,4],leaf[:,7])
plt.show()
```

Motivation
00000000

Principal component analysis
000000●0000000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

# Principal Component Analysis

Practical examples with Python



The Leaf data set

Motivation
00000000

Principal component analysis
0000000●00000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

# Principal Component Analysis
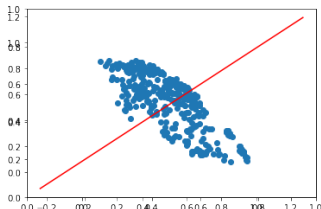Practical examples with Python
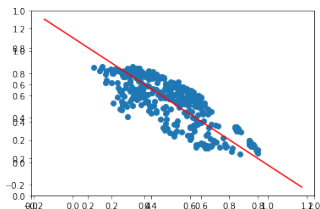
### Some questions

- How can we summarize properly using only one variable the information in this data set?
- Which variable allows to separate in the best possible way the data?
- Can we find an orientation along which the variance of the data is much higher ?

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○●○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis

Practical examples with Python

Several possibilities....

...the axis of the figure on the left seems to be the best one!

Motivation
00000000

Principal component analysis
0000000000●0000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
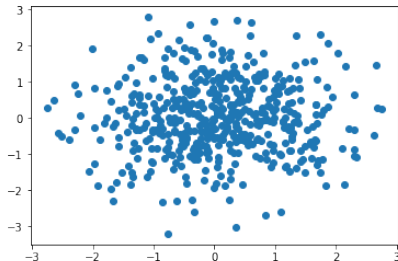0000000000

# Principal Component Analysis
Practical examples with Python

We try with a synthetic dataset!

```python
rndn = np.random.randn(500,2)
fig, ax = plt.subplots()
ax.scatter(rndn[:,0],rndn[:,1])
plt.show()
```

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○●○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis
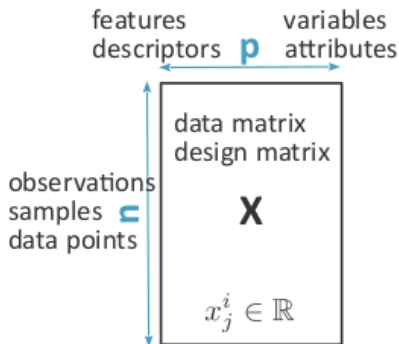Practical examples with Python



Not always possible to find an axis separating properly the data!

# Principal Component Analysis
## Principle

How summarize information in a large dataset (*n* large) living in an high dimensional space (*p* large)?

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○●○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis
Principle

## Principal Component Analysis : how does it work?

- Denote $X$ the $n \times p$ data matrix or design matrix
- We want to find specific directions maximizing the variability of the data which is summarized in the covariance matrix $X^T X$
- The $k$ dimensional space $H$ that we are looking for is generated by the $k$ eigenvectors $u_i$, $i = 1, \cdots, k$ associated to the $k$ largest eigenvalues $\lambda_i$ of the matrix $X^T X$

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○●○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis
Principle

## Principal component analysis : how does it work?

- The eigenvectors $u_1, \cdots, u_k$ are called the $k$ principal components.

- The eigenvalues $\lambda_1 \geq \cdots \geq \lambda_k$ are the explained variance ratio corresponding to each principal component

- By definition, the $k$ top principal components contain higher variance from the data.

- PCA can then be used as filtering, by selecting only the top significant PCs

Motivation
00000000

Principal component analysis
0000000000000●00000

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

# Principal Component Analysis
Principle

### Principal component analysis : how does it work?

Several possible choices for the matrix $X$ :

- General PCA : the raw data matrix $X = R$
- Centered PCA: the centered data matrix. the matrix $X^T X$ is then the matrix of empirical covariances
- Normed PCA : the normed and centered data matrix. The matrix $X^T X$ is then the matrix of empirical correlations

Motivation
00000000

Principal component analysis
0000000000000●0000

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

# Principal Component Analysis
Principle

### Principal component analysis : how does it work?

- In general $n \gg p$ (number of observations $\gg$ number of initial variables)
- It is the reason why we deal with the matrix $X^T X$ with dimension $p \times p$ rather than $XX^T$ with dimension $n \times n$
- These two analysis can both be deduced from the Singular Value Decomposition of $X$

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○●○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis

Practical examples with Python

We import the data from the website :
`https://archive.ics.uci.edu/ml/machine-learning-databases/`
`iris/iris.data`

|  | sepal length | sepal width | petal length | petal width | target |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |
| ... | ... | ... | ... | ... | ... |
| 145 | 6.7 | 3.0 | 5.2 | 2.3 | Iris-virginica |
| 146 | 6.3 | 2.5 | 5.0 | 1.9 | Iris-virginica |
| 147 | 6.5 | 3.0 | 5.2 | 2.0 | Iris-virginica |
| 148 | 6.2 | 3.4 | 5.4 | 2.3 | Iris-virginica |
| 149 | 5.9 | 3.0 | 5.1 | 1.8 | Iris-virginica |

150 rows × 5 columns

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○●○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○○

# Principal Component Analysis
Practical examples with Python

Application of a two components PCA and visualization

Motivation
00000000

Principal component analysis
000000000000000000●0

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

# Principal Component Analysis
Link with different species?

Motivation
00000000

Principal component analysis
00000000000000000000●

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000000

## Principal Component Analysis
Pro and cons of PCA

### Main advantages of PCA

- Simple to implement, no tuning
- Highly interpretable. We can find decide on how much variance to preserve using eigenvalues.

### Main drawbacks of PCA

- It is a global transform which may not preserve local structure (clusters)
- It is sensitive to outliers

An alternative : t-distributed stochastic neighbourhood embedding (t-SNE)

Motivation
OOOOOOOO

Principal component analysis
OOOOOOOOOOOOOOOOOOOO

t-distributed stochastic neighbourhood embedding (t-SNE)
●OOOOOOOOO

# t-distributed stochastic neighbourhood embedding (t-SNE)

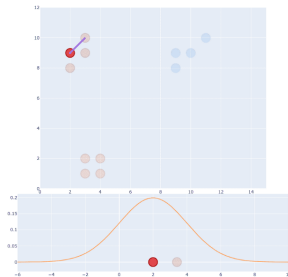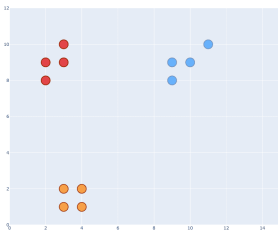More on t-SNE : https://lvdmaaten.github.io/tsne/

### Principle of SNE (Hinton 2002)

- Core idea : define an embedding from a high dimensional space to a low dimensional one, so that neighborhood identities are preserved.
- Similarity in the high dimensional space of two observations $x_j$ to $x_i$ is conditional probability $p_{j|i}$ that $x_i$ would pick $x_j$ as its nearest neighbor
- This conditional probability $p_{\cdot|i}$ is a Gaussian and depends on the (known) relative positions of the observations in the original space

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○●○○○○○○○○○

# t-distributed stochastic neighbourhood embedding (t-SNE)

More on t-SNE : `https://lvdmaaten.github.io/tsne/`

Definition of a probability $p_{.|i}$ associated to each observation $x_i$

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○●○○○○○○○

# t-distributed stochastic neighbourhood embedding (t-SNE)

More on t-SNE : https://lvdmaaten.github.io/tsne/

### The SNE mapping (Hinton 2002)?

- $y_i$ : counterpart of $x_i$ in the low dimensional space, $q_{j|i}$ : conditional probability that $y_i$ would pick $y_j$ as its nearest neighbor

- If the map points $y_i$ and $y_j$ correctly model the similarity between the high-dimensional data points $x_i$ and $x_j$, the conditional probabilities $p_{j|i}$ and $q_{j|i}$ will be equal, or at least close to each other

- Goal of SNE : find a low-dimensional data representation that minimizes the mismatch between $p_{j|i}$ and $q_{j|i}$ minimizing

$$\sum_i KL(p_{\cdot|i}, q_{\cdot|i}) \text{ (KL Kullback Divergence)}$$

Motivation
00000000

Principal component analysis
000000000000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
000●00000

# t-distributed stochastic neighbourhood embedding (t-SNE)

More on t-SNE : https://lvdmaaten.github.io/tsne/

- Leads to an optimization problem quite difficult to solve (in particular non symmetric)
- Introduction of t-SNE. Two main differences with SNE
  - uses a symmetrized version of the SNE cost function
  - uses a Student-t distribution rather than a Gaussian to compute the similarity between two points in the low-dimensional space
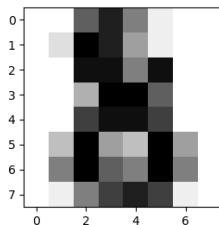
Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○●○○○○○

# t-distributed stochastic neighbourhood embedding (t-SNE)

Practical examples with Python : comparison of PCA and t-SNE on the dataset digits
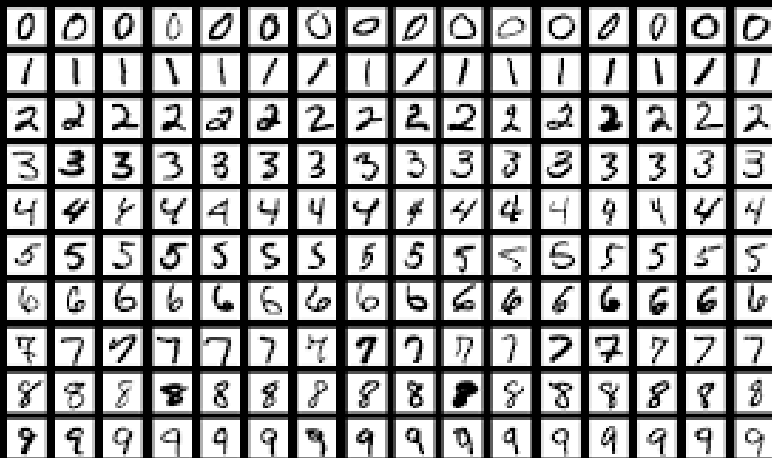
## The Digit Dataset

This dataset is made up of 1797 8x8 images. Each image, like the one shown below, is of a hand-written digit. In order to utilize an 8x8 figure like this, we'd have to first transform it into a feature vector with length 64.

See here for more information about this dataset.

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○●○○○○

# t-distributed stochastic neighbourhood embedding (t-SNE)

Practical examples with Python : comparison of PCA and t-SNE on the dataset digits

Motivation
00000000

Principal component analysis
000000000000000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
0000000●000

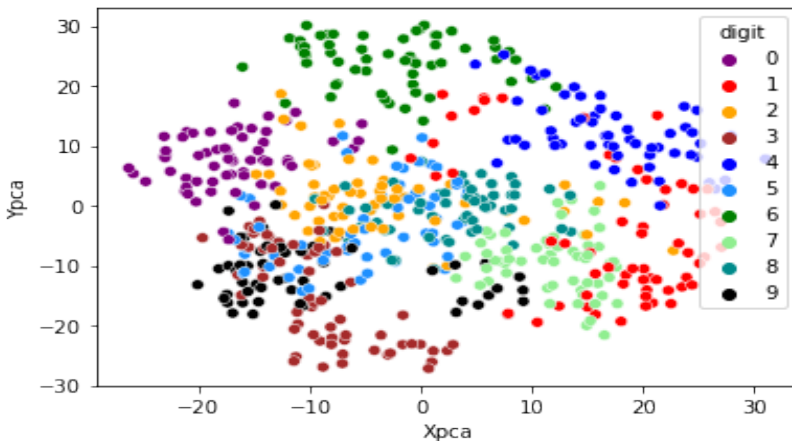# t-distributed stochastic neighbourhood embedding (t-SNE)

Practical examples with Python : comparison of PCA and t-SNE on the dataset digits

We now compare

- Two dimensional PCA
- t-SNE embedding in a two-dimensional space

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○●○○

# t-distributed stochastic neighbourhood embedding (t-SNE)

Practical examples with Python : comparison of PCA and t-SNE on the dataset digits



The PCA answer

Motivation
○○○○○○○○

Principal component analysis
○○○○○○○○○○○○○○○○○○○○○

t-distributed stochastic neighbourhood embedding (t-SNE)
○○○○○○○○○●○

# t-distributed stochastic neighbourhood embedding (t-SNE)

Practical examples with Python : comparison of PCA and t-SNE on the dataset digits



The t-SNE answer

Motivation
00000000

Principal component analysis
0000000000000000000

t-distributed stochastic neighbourhood embedding (t-SNE)
00000000●

# t-distributed stochastic neighbourhood embedding (t-SNE)
Pro and cons of t-SNE

### Main advantages of t-SNE

- It tries to preserve the local structure(cluster) of data.
- It is one of the best dimensionality reduction technique
- It can handle outliers.

### Main drawbacks of t-SNE

- It is not deterministic and iterative so each time it runs, it could produce a different result.
- It is long to run compared to PCA
- It involves hyperparameters to tune.