

Advanced Machine Learning

Lecture 8: Dimension Reduction and Data Visualization

Nora Ouzir : nora.ouzir@centralesupelec.fr

Lucca Guardiola : lucca.guardiola@centralesupelec.fr

Oct. - Nov. 2020



CentraleSupélec

Content

1. Reminders on ML
2. Robust regression
3. Hierarchical clustering
4. Classification and supervised learning
5. Non-negative matrix factorization
6. Mixture models fitting
7. Model order selection
8. Dimension reduction and data visualization

Today's Lecture

1. Introduction
2. Reminders on PCA
3. Stochastic Neighbor Embedding

Today's lecture

1. Introduction
2. Reminders on PCA
3. Stochastic Neighbor Embedding

Why should we reduce the dimension?

Huge number of dimensions causes problems....

- ▶ Data becomes very sparse, some algorithms become meaningless (e.g., density based clustering)
- ▶ Complexity usually depends on the dimensionality
- ▶ Eliminate non-relevant and/or noisy features
- ▶ Interpretation
- ▶ Visualize structure
- ▶ Storage, processing, ...

Data can be described with fewer dimensions, without losing much of its meaning

Today's lecture

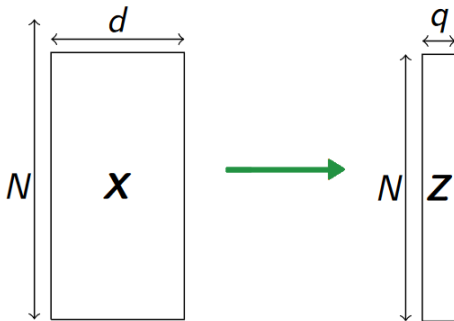
1. Introduction
2. Reminders on PCA
3. Stochastic Neighbor Embedding

PCA: principle

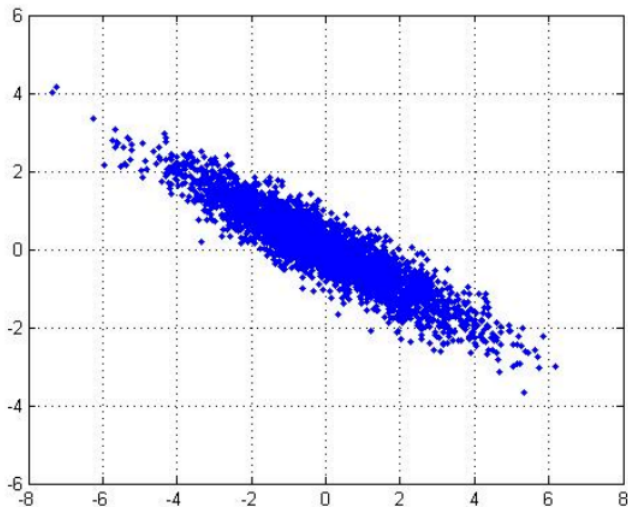
Linear orthogonal projection $\mathbf{P} \in \mathbb{R}^{d \times q}$ of the data $\mathbf{X} \in \mathbb{R}^{N, d}$ into a subspace of dimension $q < d$ s.t the variance of the projected data is maximized

$$\mathbf{X} = \mathbf{Z}\mathbf{P}^T + \text{noise}$$

where $\mathbf{Z} \in \mathbb{R}^{N \times q}$



Example: 2D Gaussian



PCA: Two equivalent views

1. **Minimization** of the reconstruction error between \mathbf{x} and $\mathbf{P}^T \mathbf{P} \mathbf{x}$

$$\min_{\mathbf{P} \in \mathbb{R}^{N \times q}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{P} \mathbf{P}^T \mathbf{x}_i\|^2$$

2. **Maximization** of variance: q largest eigenvalues of Σ

$$\max_{\mathbf{p} \in \mathbb{R}^d} \|\mathbf{X} \mathbf{p}\|_2^2 \quad \text{with} \quad \|\mathbf{p}\|_2^2 = 1 \quad \text{and} \quad \mathbf{Z} = \mathbf{X} \mathbf{p}$$

Classic matrix computation tools

- ▶ Eigenvalues of the covariance matrix $\Sigma \rightarrow$ capture variance direction and scale: $\Sigma = \mathbf{U} \Lambda \mathbf{U}^T$
- ▶ Singular Value Decomposition: $\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^T$

These are equivalent views \rightarrow

Minimization of error \rightsquigarrow Maximization of variance

$$\begin{aligned}
J(\mathbf{P}) &= \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \sum_{i=1}^N (\mathbf{x}_i - \mathbf{P}\mathbf{P}^\top \mathbf{x}_i)^\top (\mathbf{x}_i - \mathbf{P}\mathbf{P}^\top \mathbf{x}_i) \\
&= \sum_{i=1}^N (\mathbf{x}_i^\top \mathbf{x}_i - 2\mathbf{x}_i^\top \mathbf{P}\mathbf{P}^\top \mathbf{x}_i + \mathbf{x}_i^\top \mathbf{P}\mathbf{P}^\top \mathbf{P}\mathbf{P}^\top \mathbf{x}_i) \\
&= \sum_{i=1}^N \mathbf{x}_i^\top \mathbf{x}_i - \sum_{i=1}^N \mathbf{x}_i^\top \mathbf{P}\mathbf{P}^\top \mathbf{x}_i = \sum_{i=1}^N \mathbf{x}_i^\top \mathbf{x}_i - \sum_{i=1}^N \mathbf{z}_i^\top \mathbf{z}_i \\
&= \text{trace} \left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i \mathbf{z}_i^\top \right) = \text{trace} \left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top - \sum_{i=1}^N \mathbf{P}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{P} \right) \\
J(\mathbf{P}) &= \text{trace}(\mathbf{X}^\top \mathbf{X}) - \text{trace}(\mathbf{P}^\top \mathbf{X}^\top \mathbf{X} \mathbf{P})
\end{aligned}$$

Minimize $J(\mathbf{P}) \iff$ maximize variance of the projection wrt \mathbf{P}

Reminders on PCA

Whiteboard

PCA: Solution

The solution to the PCA problem is the matrix \mathbf{P} whose columns are the q eigenvectors of the covariance matrix Σ corresponding to the q largest eigenvalues

- ▶ → N.B: The solution can be obtained with SVD of the data matrix

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

Choice of q

Choose the q largest components that explain e.g., 90% of the variance in the data

$$\text{total variance} \rightarrow \text{trace}(\Lambda) = \lambda_1 + \dots + \lambda_d$$

The first q components account for: $\frac{\lambda_1 + \dots + \lambda_q}{\lambda_1 + \dots + \lambda_d}$

PCA: Algorithm

Algorithm

1. Normalize the data $\mathbf{X} \in \mathbb{R}^{N \times d}$: $\mathbf{x}_{ij} = \frac{\mathbf{x}_{ij} - \mu_j}{\sigma_j}, j = 1, \dots, d$ and $i = 1, \dots, N$
2. Compute the correlation matrix $\Sigma = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$
3. Find the eigenvalue decomposition of Σ : $\{\mathbf{p}_i \in \mathbb{R}^d, \lambda_j \in \mathbb{R}\}$
4. Order the eigenvalues λ_j by decreasing order
5. The projection matrix is

$$\mathbf{P} = (\mathbf{p}_1, \dots, \mathbf{p}_q) \in \mathbb{R}^{d \times q}$$

where $\{\mathbf{p}_1, \dots, \mathbf{p}_q\}$ are the eigenvectors associated with the q largest eigenvalues

PCA: Summary

- ▶ Linear dimension reduction method
- ▶ PCA projects the data onto a subspace which maximizes the projected variance, or equivalently, minimizes the reconstruction error.
- ▶ The optimal subspace is given by the top eigenvectors of the empirical covariance matrix.
- ▶ PCA gives a set of decorrelated features.

Extensions

- ▶ Kernel PCA: non-linear projection
- ▶ Probabilistic PCA
 - Bypass covariance matrix computation, missing data, mixtures, ...
 - Bishop 2006 Chapter 12.2

Example: Non-linear subspace

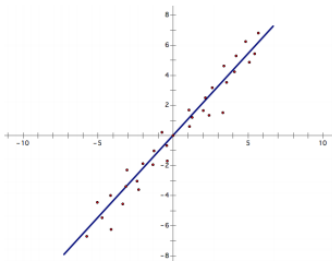


Figure 1. Data lying near a linear subspace

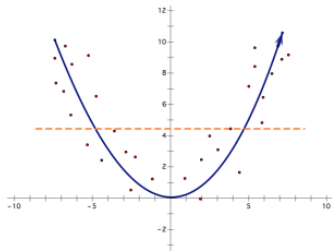


Figure 2. Data lying near a parabola

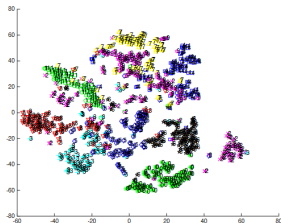
Today's lecture

1. Introduction
2. Reminders on PCA
3. Stochastic Neighbor Embedding

MNIST dataset

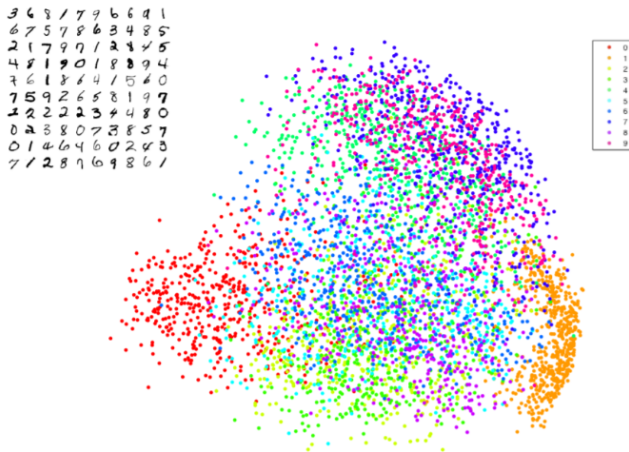


$d = 784$



$q = 2$

MNIST dataset: PCA results



What can you say about the structure of this data?

SNE: Principles

Key Idea: Treat pairwise distances as probabilities \rightarrow stochastic
Low distance \Rightarrow High probability

Principles

- ▶ Small distance does not mean proximity on manifold
- ▶ *Encode* high dimensional neighborhood information as a distribution
- ▶ Picking two similar points is much more probable than picking two points which are well far apart
- ▶ Preserve **local** structure
 - \rightarrow Low-dim and original neighborhoods should be the same

Whiteboard

SNE: Key steps

1. Transform pairwise distances $\text{dist}(\mathbf{x}_i, \mathbf{x}_j)$ into probabilities $P_X(\mathbf{x}_i|\mathbf{x}_j)$
2. Do the same for the projection space $\text{dist}(\mathbf{y}_i, \mathbf{y}_j) \rightarrow P_Y(\mathbf{y}_i|\mathbf{y}_j)$
3. Minimize the distance between the distributions P_X and $P_Y \rightarrow$ Optimization problem

SNE

SNE: Key steps

1. Transform pairwise distances $\text{dist}(\mathbf{x}_i, \mathbf{x}_j)$ into probabilities $P_X(\mathbf{x}_i|\mathbf{x}_j)$
2. Do the same for the projection space $\text{dist}(\mathbf{y}_i, \mathbf{y}_j) \rightarrow P_Y(\mathbf{y}_i|\mathbf{y}_j)$
3. Minimize the distance between the distributions P_X and $P_Y \rightarrow$ Optimization problem

SNE

- ▶ Compare distributions \rightarrow Kullback-Leibler (KL) divergence

SNE: Key steps

1. Transform pairwise distances $\text{dist}(\mathbf{x}_i, \mathbf{x}_j)$ into probabilities $P_X(\mathbf{x}_i|\mathbf{x}_j)$
2. Do the same for the projection space $\text{dist}(\mathbf{y}_i, \mathbf{y}_j) \rightarrow P_Y(\mathbf{y}_i|\mathbf{y}_j)$
3. Minimize the distance between the distributions P_X and $P_Y \rightarrow$ Optimization problem

SNE

- ▶ Compare distributions \rightarrow Kullback-Leibler (KL) divergence
- ▶ Model P_X and $P_Y \rightarrow$ Gaussian distributions

SNE: Key steps

1. Transform pairwise distances $\text{dist}(\mathbf{x}_i, \mathbf{x}_j)$ into probabilities $P_X(\mathbf{x}_i|\mathbf{x}_j)$
2. Do the same for the projection space $\text{dist}(\mathbf{y}_i, \mathbf{y}_j) \rightarrow P_Y(\mathbf{y}_i|\mathbf{y}_j)$
3. Minimize the distance between the distributions P_X and $P_Y \rightarrow$ Optimization problem

SNE

- ▶ Compare distributions \rightarrow Kullback-Leibler (KL) divergence
- ▶ Model P_X and $P_Y \rightarrow$ Gaussian distributions
- ▶ Different distribution for each sample! The size of the local neighborhood changes

Neighborhood Distributions: Gaussian case

Neighborhood Distributions: Gaussian case

Goal: model the probability that a point \mathbf{x}_i chooses \mathbf{x}_j as its neighbor

- ▶ Input space: $P_X(\mathbf{x}_j|\mathbf{x}_i) = p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)}$
- ▶ Output space: $P_Y(\mathbf{y}_j|\mathbf{y}_i) = q_{j|i} = \frac{\exp(-\|\mathbf{y}_i - \mathbf{y}_j\|^2)}{\sum_{k \neq i} \exp(-\|\mathbf{y}_i - \mathbf{y}_k\|^2)}$
- ▶ $p_{i|i} = q_{i|i} = 0$
- ▶ Symmetrization $p_{ji} = \frac{1}{2N}(p_{j|i} + p_{i|j})$

Neighborhood Distributions: Gaussian case

Goal: model the probability that a point \mathbf{x}_i chooses \mathbf{x}_j as its neighbor

- ▶ Input space: $P_X(\mathbf{x}_j|\mathbf{x}_i) = p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)}$
- ▶ Output space: $P_Y(\mathbf{y}_j|\mathbf{y}_i) = q_{j|i} = \frac{\exp(-\|\mathbf{y}_i - \mathbf{y}_j\|^2)}{\sum_{k \neq i} \exp(-\|\mathbf{y}_i - \mathbf{y}_k\|^2)}$
- ▶ $p_{i|i} = q_{i|i} = 0$
- ▶ Symmetrization $p_{ji} = \frac{1}{2N}(p_{j|i} + p_{i|j})$

How to choose σ_i ?

- ▶ Different value for each point i
- ▶ Very **low** $\sigma_i \rightarrow$ all the probability is in the nearest neighborhood
- ▶ Very **high** $\sigma_i \rightarrow$
- ▶ In practice \rightarrow **perplexity**

Perplexity

Perplexity is a smooth measure of the number of neighbors

$$\text{Perp}(P_i) = 2^{-\sum_j p_{ji} \log_2(p_{ji})}$$

- ▶ Uniform P over k elements $\rightarrow \text{Perp}(P) = k$
- ▶ High/low perplexity \rightarrow Large/small σ_i
- ▶ In practice: values between 5 and 50



Perplexity should be smaller than N !

SNE: Optimization Problem

Minimize the KL divergence between p_{ji} and q_{ji}

$$\min_{y_1, \dots, y_N} \mathcal{C}(X, Y) = \sum_{i,j} KL(P_{ji}|Q_{ji}) = \sum_{i,j} p_{ji} \log \frac{p_{ji}}{q_{ji}}$$

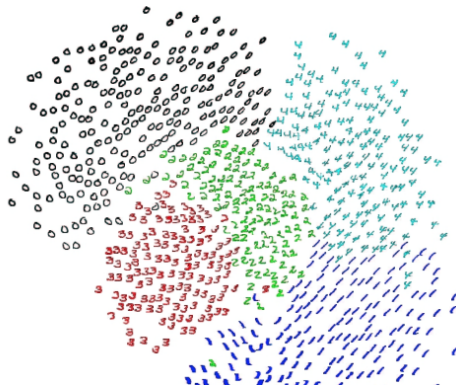
→ One can show that (symmetric case)

$$\frac{\partial \mathcal{C}}{\partial \mathbf{y}_i} = \sum_j (p_{ji} - q_{ji})(\mathbf{y}_i - \mathbf{y}_j)$$

→ Non-convex: no convergence guarantees, multiple restarts

Solved with numerical methods

SNE: MNIST results



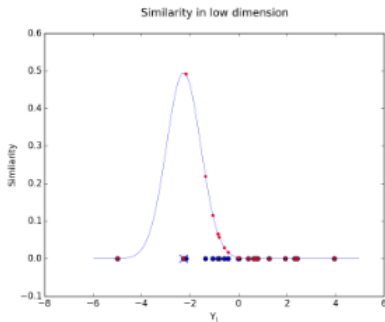
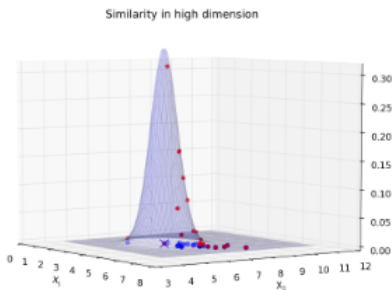
But...

- ▶ Impossible to preserve the whole structure due to the distortions introduced by the dimensionality reduction
- ▶ **Crowding problem!**

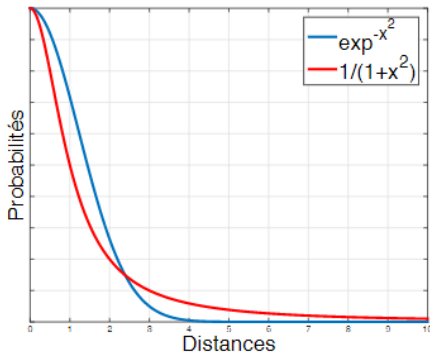
SNE: Limitation

Crowding problem

- ▶ We don't have enough room to accommodate all neighbors
- ▶ Example: embed a circular neighborhood from 2D to 1D



t-SNE: dealing with the crowding problem



Heavy-tailed q distribution

- ▶ Local neighbors are **attracted** to a slightly closer position
- ▶ Non-local neighbors might get a little **repulsed**
- ▶ Non-neighbors are potentially pushed away

t-SNE principles

Key Idea: model the output space using a heavy-tailed distribution

→ Student's t-distribution $\nu = 1$

- ▶ Input space: $P_X(\mathbf{x}_j|\mathbf{x}_i) = p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)}$
- ▶ Output space: $P_Y(\mathbf{y}_j|\mathbf{y}_i) = q_{j|i} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|\mathbf{y}_i - \mathbf{y}_k\|^2)^{-1}}$
- ▶ $p_{i|i} = q_{i|i} = 0$
- ▶ Symmetrization $p_{ji} = \frac{1}{2N}(p_{j|i} + p_{i|j})$

Attraction/repulsion

- ▶ large probability → $d_Y < d_X$
- ▶ low probability → $d_Y > d_X$

t-SNE: Optimization Problem

Minimize the KL divergence between p_{ji} and q_{ji}

$$\min_{y_1, \dots, y_N} C(X, Y) = \sum_{i,j} KL(P_{ji} | Q_{ji}) = \sum_{i,j} p_{ji} \log \frac{p_{ji}}{q_{ji}}$$

→ Gradient

$$\frac{\partial C}{\partial \mathbf{y}_i} = \sum_j (p_{ji} - q_{ji})(\mathbf{y}_i - \mathbf{y}_j)(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}$$

→ Non-convex: no convergence guarantees, multiple restarts

Solved with gradient descent

t-SNE: Algorithm

Algorithm 1: Simple version of t-Distributed Stochastic Neighbor Embedding.

Data: data set $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$,

cost function parameters: perplexity $Perp$,

optimization parameters: number of iterations T , learning rate η , momentum $\alpha(t)$.

Result: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, \dots, y_n\}$.

begin

 compute pairwise affinities $p_{j|i}$ with perplexity $Perp$ (using Equation 1)

 set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

 sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, \dots, y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$

for $t=1$ **to** T **do**

 compute low-dimensional affinities q_{ij} (using Equation 4)

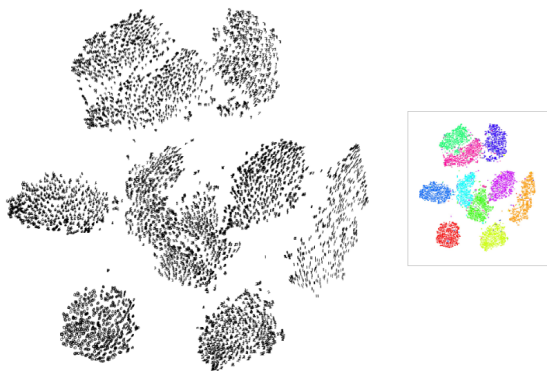
 compute gradient $\frac{\partial C}{\partial \mathcal{Y}}$ (using Equation 5)

 set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\partial C}{\partial \mathcal{Y}} + \alpha(t) (\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)})$

end

end

t-SNE: MNIST results



But...

- ▶ Different perplexity → different results
- ▶ Different results with different runs

PCA vs t-SNE

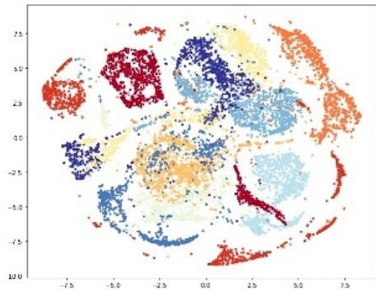
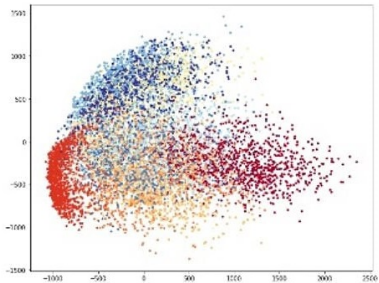
PCA

- ▶ Finds **global** structure
- ▶ Low-dim subspace
 - Inconsistencies: Far away points can become nearest neighbors
 - Doesn't always reveal hidden/complex structure (linear)
- + Relatively fast
- + Features can be traced back

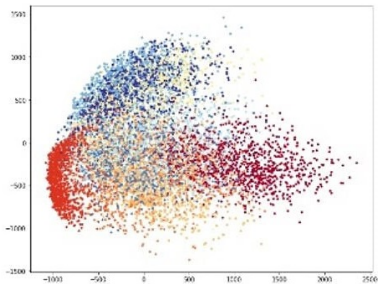
t-SNE

- ▶ Preserves **local** structure
- ▶ Low and high-dim neighborhoods should be the same
- ▶ Used for visualization
- + Reveals hidden structure
 - Relatively slow ($O(N^2)$)
 - Meaning of the features is lost
 - Cannot project a new point onto an already computed map (+ variations exist)

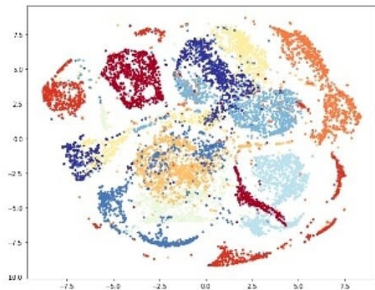
Example



Example



PCA



t-SNE