

Data Mining

Visualization of Multidimensional Data

2/12/2021

NOVA-IMS

Fernando Lucas Bação
bacao@isegi.unl.pt
<http://www.isegi.unl.pt/fbacao>

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

1

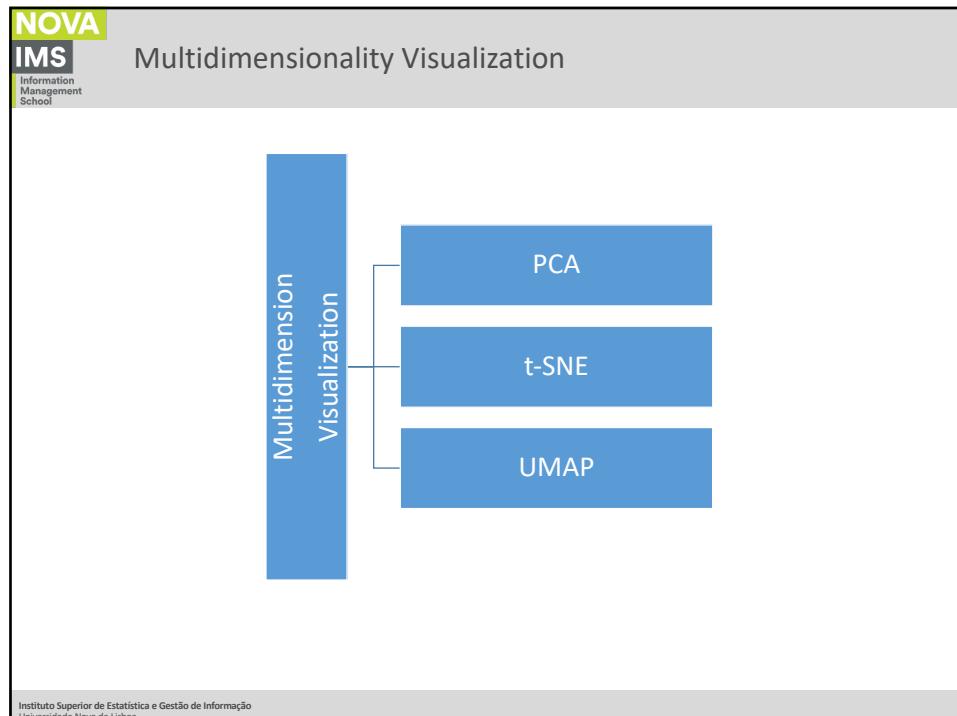


Multidimensionality Visualization

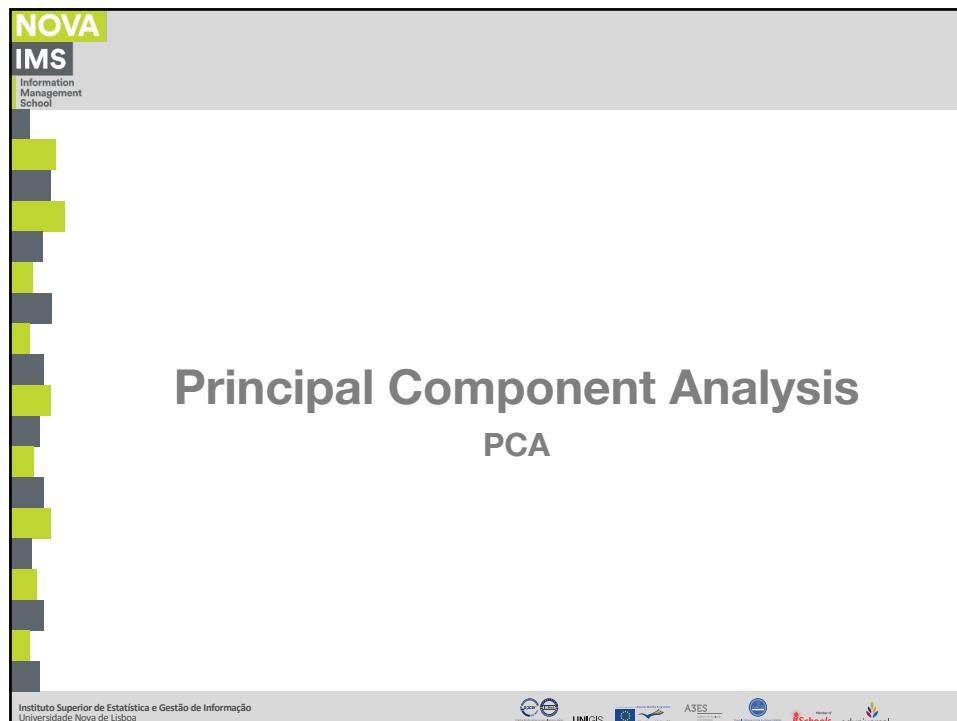
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

AACSB Accredited
UNICIS
A3ES
EQUIS Accredited
World's Best Schools
eduniversal

2



3



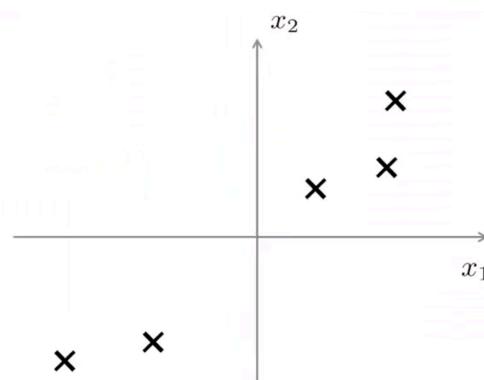
4

- **Size Reduction of the Input Space:**

- Principal Component Analysis
- A procedure that uses an **orthogonal transformation** to convert a set of observations of possibly correlated variables into a set of values of **linearly uncorrelated variables called principal components**.
- The number of principal components is **equal to the number of original variables**.
- This transformation is defined in such a way that the first principal component has the **largest possible variance** (that is, accounts for as much of the variability in the data as possible), and each succeeding component in turn has the highest variance.

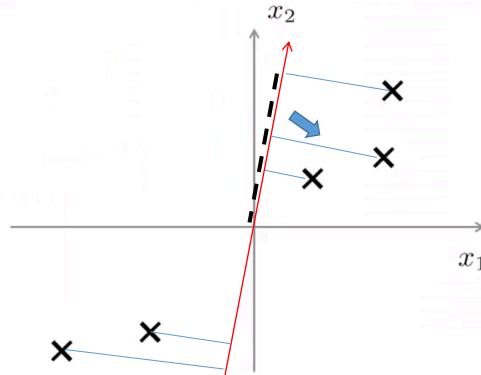
- **Size Reduction of the Input Space:**

- Principal Component Analysis



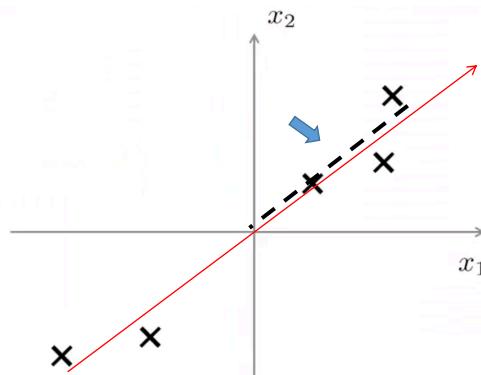
- **Size Reduction of the Input Space:**

- Principal Component Analysis (SVD)



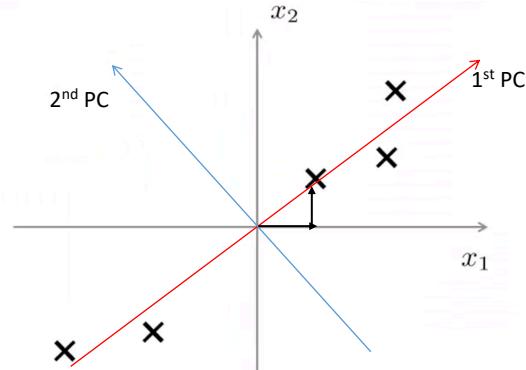
- **Size Reduction of the Input Space:**

- Principal Component Analysis (SVD)



- **Size Reduction of the Input Space:**

- Principal Component Analysis (SVD)

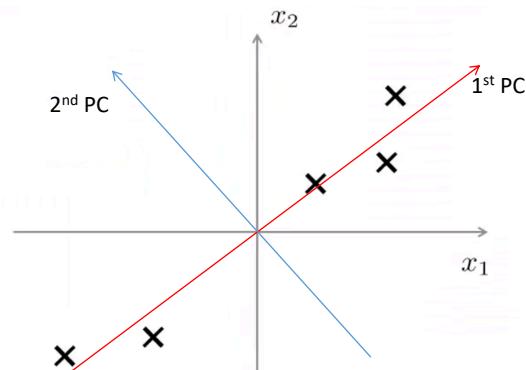


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

9

- **Size Reduction of the Input Space:**

- Principal Component Analysis (SVD)

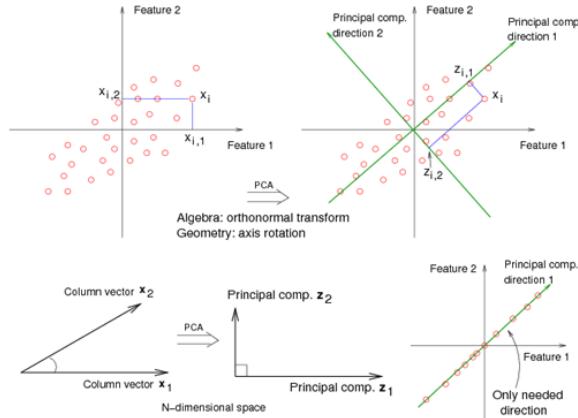


Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

10

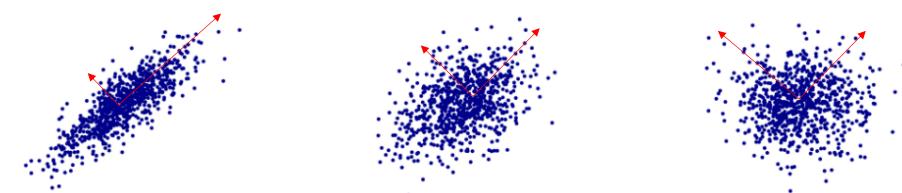
- **Size Reduction of the Input Space:**

- Principal Component Analysis



- **Size Reduction of the Input Space:**

- Principal Component Analysis



NOVA
IMS
Information Management School

t-distributed Stochastic Neighbor Embedding

t-SNE

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

EQUIS UNIGIS A3ES AACSB iSchools eduniversal

13

NOVA
IMS
Information Management School

t-SNE

- **t-SNE:**
 - t-distributed stochastic neighbor embedding (t-SNE) is a statistical method for visualizing high-dimensional data by **giving each datapoint a location in a two or three-dimensional map.**
 - It is a **nonlinear dimensionality reduction technique** well-suited for embedding high-dimensional data for visualization in a low-dimensional space of two or three dimensions.
 - Specifically, it **models each high-dimensional object** by a two- or three-dimensional point in such a way that **similar objects are modeled by nearby points** and dissimilar objects are modeled by distant points with high probability.

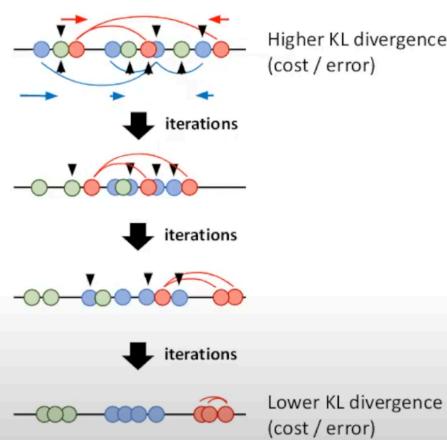
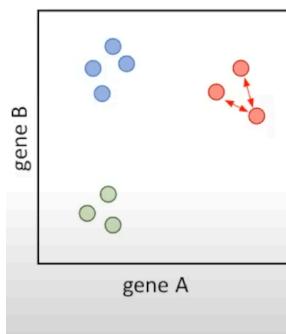
Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

14

- **t-SNE:**

- The t-SNE algorithm comprises two stages:
 - First, t-SNE constructs a **probability distribution over pairs of high-dimensional objects** in such a way that similar objects are assigned a higher probability while dissimilar points are assigned a lower probability.
 - Second, t-SNE defines a **similar probability distribution over the points in the low-dimensional map**, and it minimizes the Kullback–Leibler divergence (KL divergence) between the two distributions with respect to the locations of the points in the map.
- While the original algorithm uses the Euclidean distance between objects as the base of its similarity metric, this can be changed as appropriate.

- **t-SNE:**



NOVA
IMS
Information Management School

t-SNE

- **t-SNE:**
 - Optimizing the KL divergence is a measure of how one probability distribution is different from a second, reference probability distribution

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

17

NOVA
IMS
Information Management School

t-SNE

- **t-SNE:**
 - t-SNE has been used for **visualization in a wide range of applications**, including genomics, computer security research, natural language processing, music analysis, cancer research, bioinformatics, etc
 - While t-SNE plots often seem to display clusters, **the visual clusters can be influenced strongly by the chosen parameterization** and therefore a good understanding of the parameters for t-SNE is necessary.
 - Interactive exploration may thus be necessary to choose parameters and validate results.

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

18

**NOVA
IMS**
Information Management School

t-SNE

- **t-SNE:**

 - Perplexity is the main t-SNE parameter, “perplexity,”
 - Basically defines how to balance attention between local and global aspects of your data.
 - Perplexity is a measure for information that is defined as 2 to the power of the Shannon entropy.
 - In t-SNE, the perplexity may be viewed as a knob that sets the number of effective nearest neighbors.
 - It is comparable with the number of nearest neighbors k that is employed in many manifold learners.
 - The original paper says, “The performance of SNE is fairly robust to changes in the perplexity, and typical values are between 5 and 50.”

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

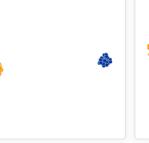
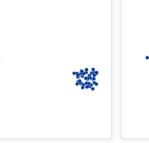
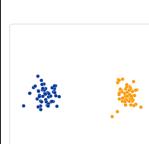
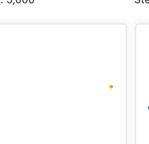
19

**NOVA
IMS**
Information Management School

t-SNE

- **t-SNE:**

 - Perplexity is the main t-SNE parameter, “perplexity,”

					
Original	Perplexity: 2 Step: 5,000	Perplexity: 5 Step: 5,000	Perplexity: 30 Step: 5,000	Perplexity: 50 Step: 5,000	Perplexity: 100 Step: 5,000
					
Original	Perplexity: 30 Step: 10	Perplexity: 30 Step: 20	Perplexity: 30 Step: 60	Perplexity: 30 Step: 120	Perplexity: 30 Step: 1,000

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

20

**NOVA
IMS**
Information Management School

t-SNE

- **t-SNE:**
 - Perplexity is the main t-SNE parameter, “perplexity,”

Perplexity	Step
Original	
2	5,000
5	5,000
30	5,000
50	5,000
100	5,000

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

21

**NOVA
IMS**
Information Management School

Uniform Manifold Approximation and Projection

UMAP

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

Accredited by:

22

- **UMAP:**

- Uniform Manifold Approximation and Projection (UMAP) is a dimension reduction technique that can be used for visualisation similarly to t-SNE.
- The algorithm is founded on three assumptions about the data
 - The data is uniformly distributed on Riemannian manifold (a topological space that locally resembles Euclidean space near each point);
 - The manifold is locally connected.
 - From these assumptions it is possible to model the manifold with a fuzzy topological structure. The embedding is found by searching for a low dimensional projection of the data that has the closest possible equivalent fuzzy topological structure.
- For a more detailed explanation see <https://youtu.be/nq6iPZVUxZU>

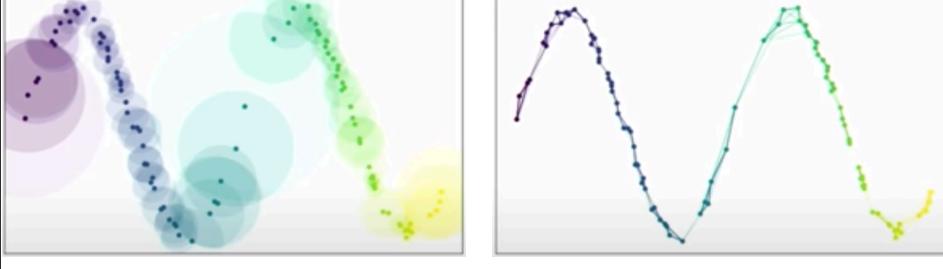
- **UMAP:**

- Uniform Manifold Approximation and Projection (UMAP) is a dimension reduction technique that can be used for visualisation similarly to t-SNE.
- Two phases in the UMAP algorithm
 - The first phase consists of constructing a fuzzy topological representation in the original space;
 - The second phase is simply optimizing (through stochastic gradient descent) the low dimensional representation to have as close a fuzzy topological representation as possible as measured by cross entropy.

NOVA
IMS
Information Management School

UMAP

- **UMAP:**
 - The first phase consists of constructing a fuzzy topological representation;



Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

25

NOVA
IMS
Information Management School

UMAP

- **UMAP:**
 - The second phase is simply optimizing (through stochastic gradient descent) the low dimensional representation to have as close a fuzzy topological representation as possible as measured by cross entropy.

Get the clumps right

$$\sum_{a \in A} \mu(a) \log \left(\frac{\mu(a)}{\nu(a)} \right) + (1 - \mu(a)) \log \left(\frac{1 - \mu(a)}{1 - \nu(a)} \right)$$

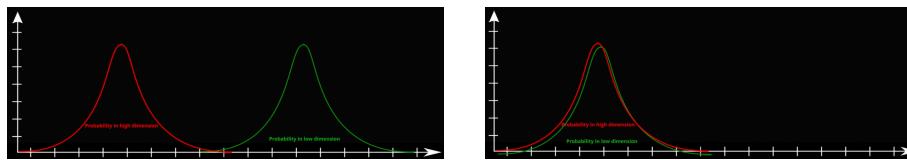
Get the gaps right

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

26

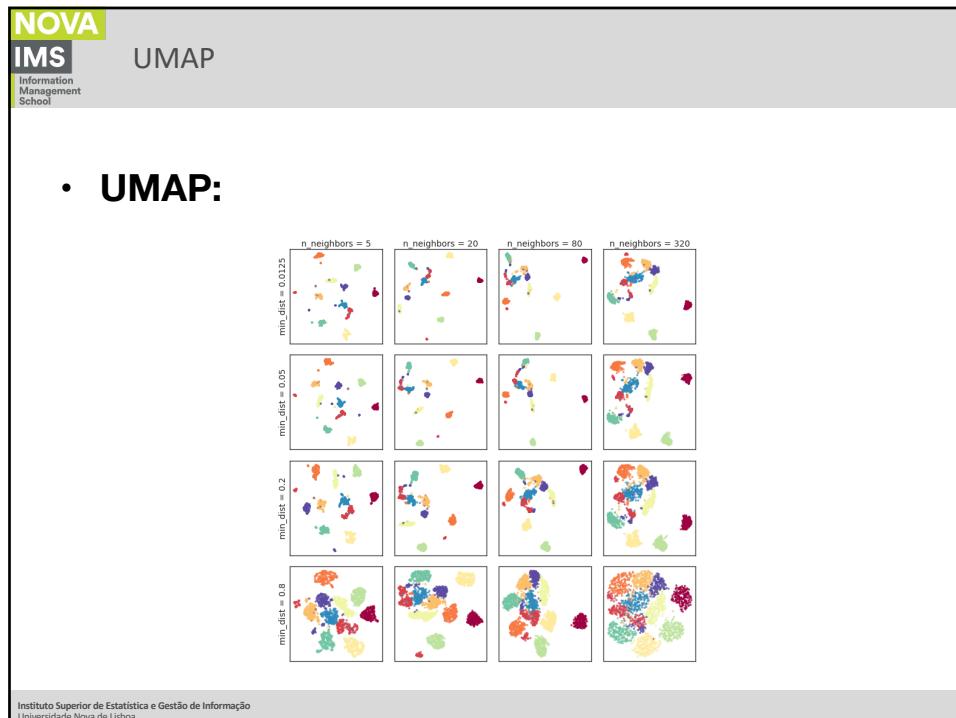
- **UMAP:**

- The second phase is simply optimizing (through stochastic gradient descent) the low dimensional representation to have as close a fuzzy topological representation as possible as measured by cross entropy.

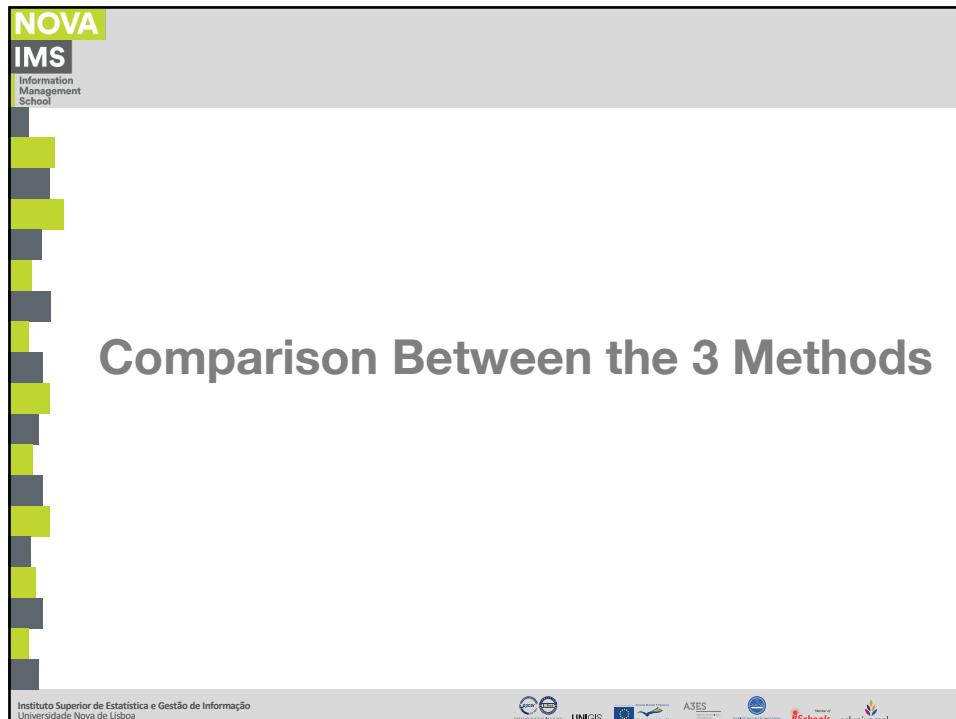


- **UMAP:**

- Uniform Manifold Approximation and Projection (UMAP) is a dimension reduction technique that can be used for visualisation similarly to t-SNE.
- Parameters
 - Number of nearest neighbors - controls how UMAP balances local versus global structure in the data;
 - Minimum distance - controls how tightly UMAP is allowed to pack points together. It provides the minimum distance apart that points are allowed to be in the low dimensional representation



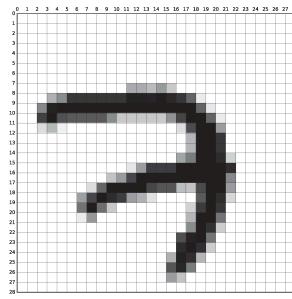
29



30

- **Mnist Digits:**

- 28x28 image (784 dimensions)
- grayscale images of handwritten single digits
- 60,000 examples, and a test set of 10,000 examples



(a) MNIST sample belonging to the digit '7'.

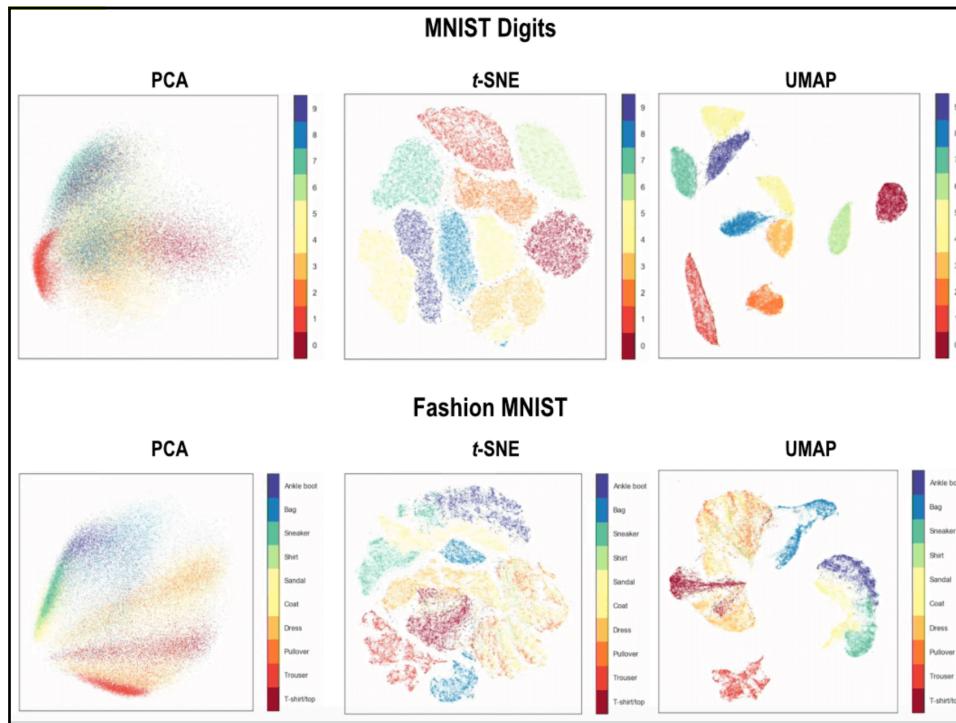


(b) 100 samples from the MNIST training set.

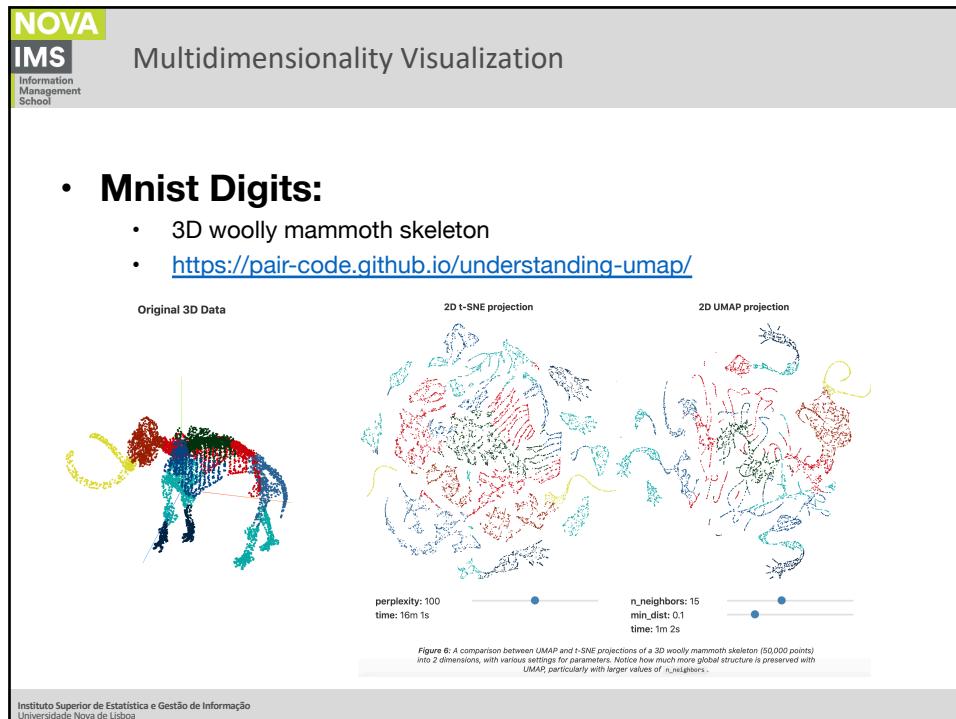
- **Mnist Fashion:**

- 28x28 image (784 dimensions)
- grayscale images of handwritten single digits
- The dataset has 60,000 images





33



34

**NOVA
IMS**
Information Management School

UMAP

- **UMAP:**

	t-SNE	UMAP
COIL20	20 seconds	7 seconds
MNIST	22 minutes	98 seconds
Fashion MNIST	15 minutes	78 seconds
GoogleNews	4.5 hours	14 minutes

Instituto Superior de Estatística e Gestão de Informação
Universidade Nova de Lisboa

35



36