

Dimensionality Reduction

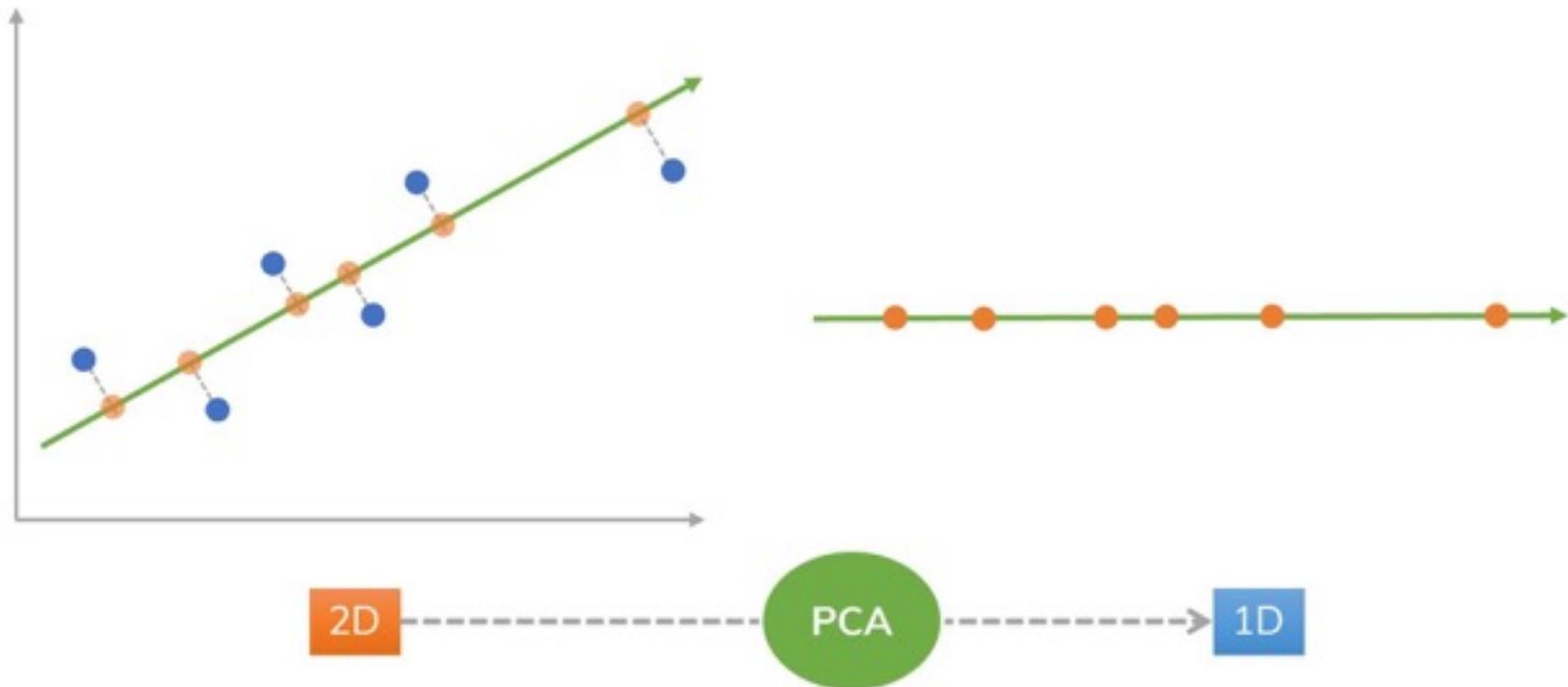


Seyed Abbas Hosseini
Sharif University of Technology

Most slides are adopted from Soleymani CE-717

Dimensionality Reduction

- **Dimensionality reduction** is the transformation of data from a high-dimensional space into a low-dimensional space so that the low-dimensional representation retains some meaningful properties of the original data.



Feature Selection vs. Extraction

- Feature Selection
 - select a subset of a given feature set
- Feature Extraction
 - A linear or nonlinear transformation on original feature set

$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \begin{bmatrix} x_{i_1} \\ \vdots \\ x_{i_{d'}} \end{bmatrix}$$

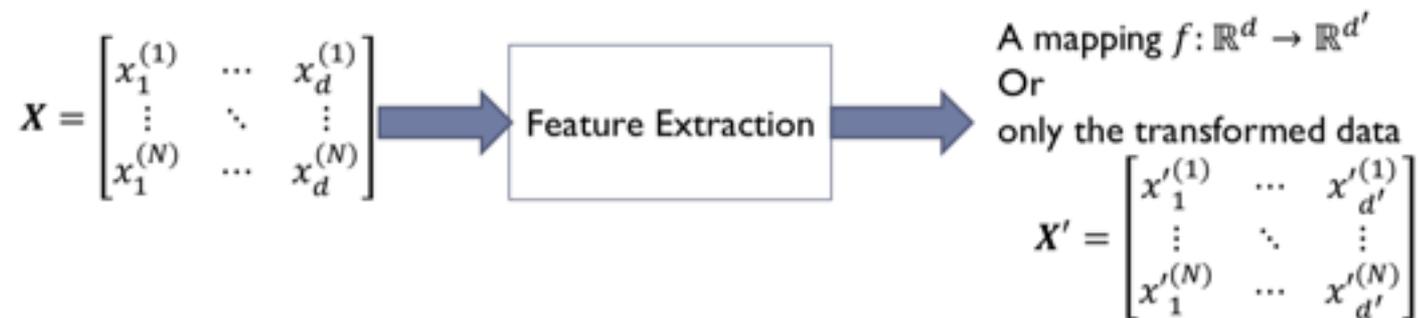
Feature
Selection
($d' < d$)

$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \begin{bmatrix} y_1 \\ \vdots \\ y_{d'} \end{bmatrix} = f \left(\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \right)$$

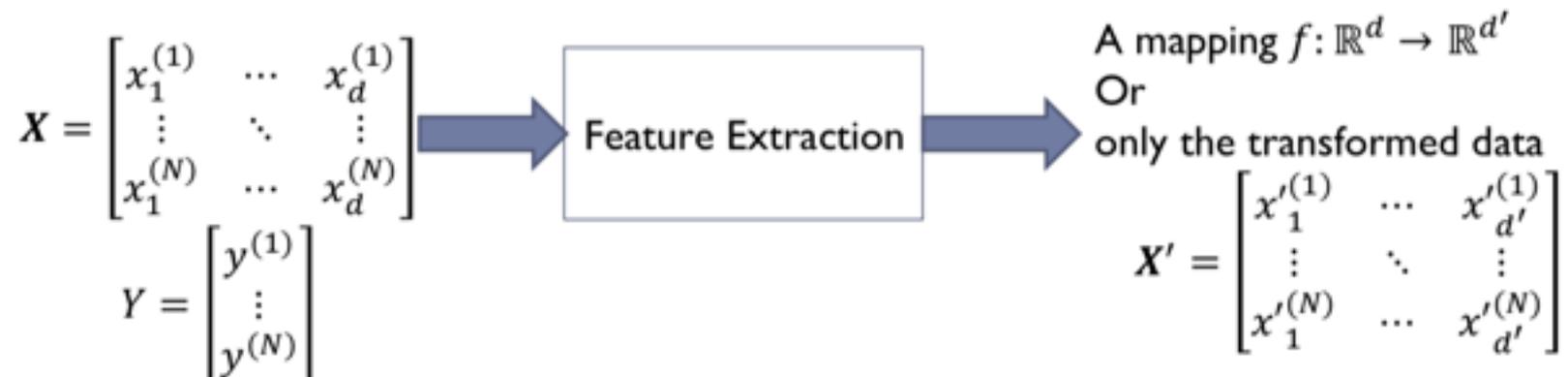
Feature
Extraction

Feature Selection vs. Extraction

- Unsupervised Feature Extraction

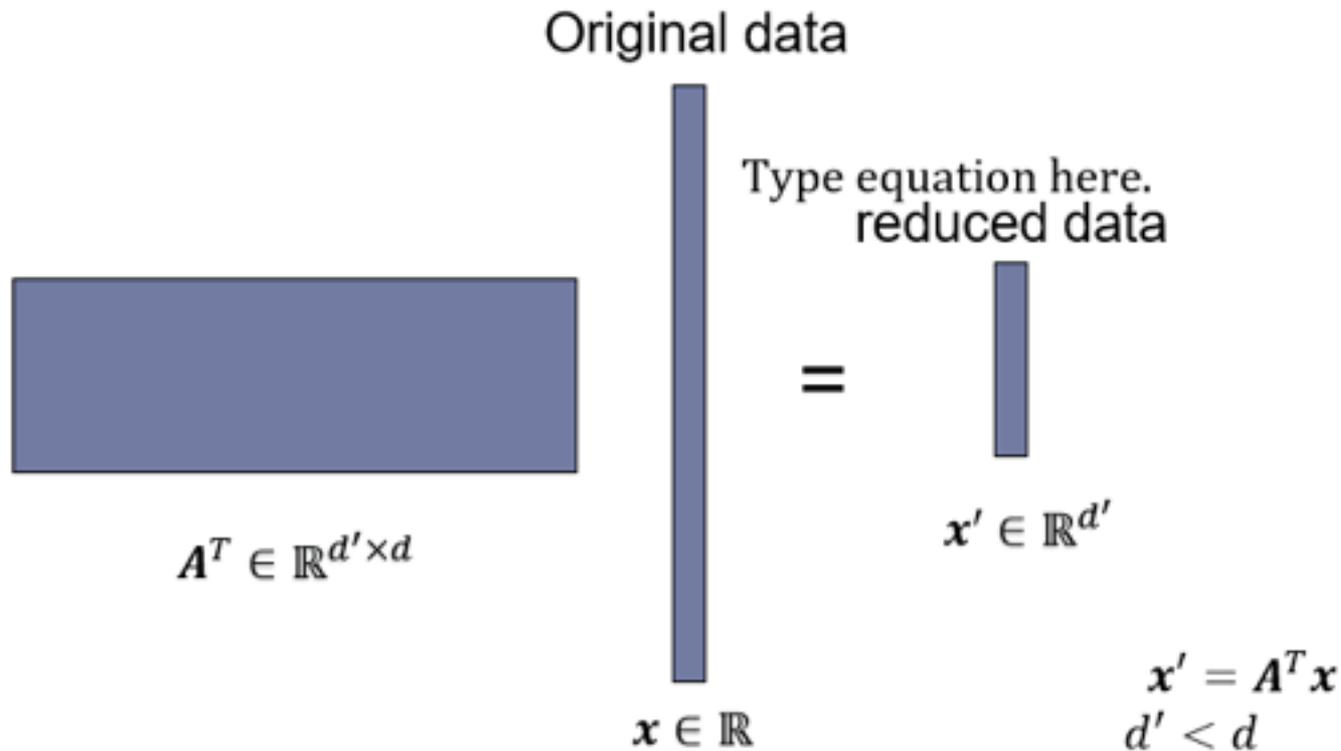


- Supervised Feature Extraction



Linear Dimensionality Reduction

- Finding a linear transformation (matrix) that directly maps data from high dimensional feature space to a low dimensional feature space



Principal Component Analysis (PCA)

- Goal:
 - ▶ Reducing the dimensionality of data by a **linear transformation** while preserving important aspects of the data

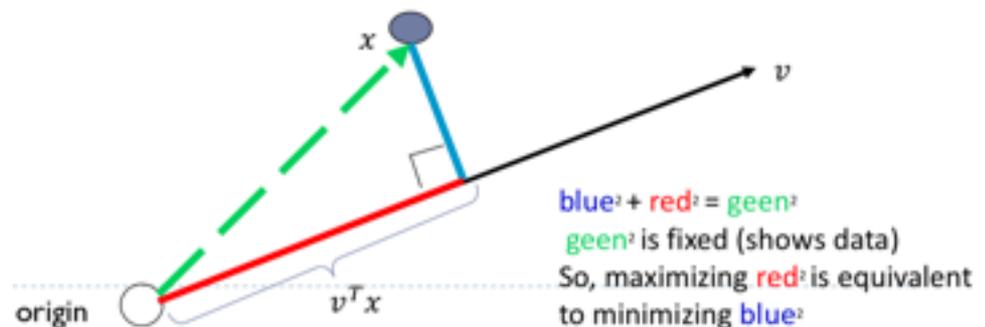
- Two equivalent views: find the direction for which

- ▶ the variation presents in the dataset is as much as possible

$$\operatorname{argmax}_v \frac{1}{N} \sum_{n=1}^N (v^T x^{(n)})^2$$

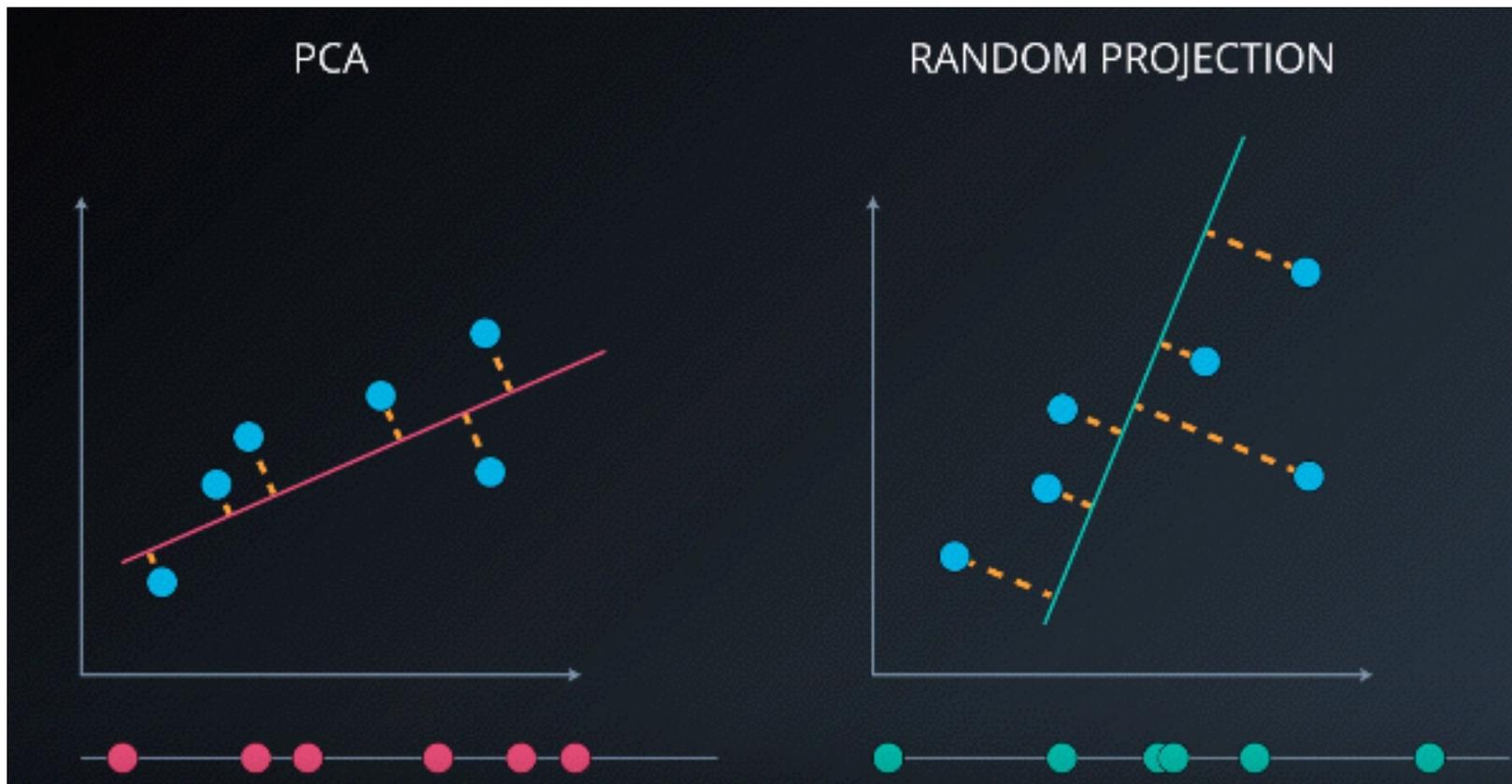
- ▶ the reconstruction error is minimized

$$\operatorname{argmin}_v \sum_{n=1}^N \|x^{(n)} - (v^T x^{(n)})v\|^2$$



PCA vs. Random Projection

PCA finds directions that preserve the information in data and original data can be reconstructed by minimum error.



Principal Component Analysis (PCA)

- Goal:
 - Reducing the dimensionality of data by a linear transformation while preserving important aspects of the data
- Two equivalent views: find the direction for which
 - the variation presents in the dataset is as much as possible
 - the reconstruction error is minimized
- Method: Mapping each data onto first few *Principal components*

Principal Components

- **Principal Components:** Orthonormal basis for the vector space of data that are ordered by the fraction of the total information (variation) in the corresponding directions
- **Claim:** PCs are the eigenvectors of the **covariance Matrix** of the data points.

$$\boldsymbol{\mu}_x = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix} = \begin{bmatrix} E(x_1) \\ \vdots \\ E(x_d) \end{bmatrix}$$

$$\boldsymbol{\Sigma} = E[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T]$$

- ML estimate of covariance matrix from $\{\mathbf{x}^{(i)}\}_{i=1}^N$ data points :

$$\mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \bar{\mathbf{x}})(\mathbf{x}^{(i)} - \bar{\mathbf{x}})^T = \frac{1}{N} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})$$

$$\tilde{\mathbf{X}} = \begin{bmatrix} \tilde{\mathbf{x}}^{(1)} \\ \vdots \\ \tilde{\mathbf{x}}^{(N)} \end{bmatrix} = \begin{bmatrix} \mathbf{x}^{(1)} - \bar{\mathbf{x}} \\ \vdots \\ \mathbf{x}^{(N)} - \bar{\mathbf{x}} \end{bmatrix} \quad \bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$$

First Principal Component

Find vector \mathbf{v}_1 that maximizes sample variance

$$\begin{aligned} & \max_{\mathbf{v}_1} \frac{1}{N} \sum_{n=1}^N (v_1^T x^{(n)} - v_1^T \bar{x})^2 \\ &= \frac{1}{N} \sum_{n=1}^N v_1^T (x^{(n)} - \bar{x})(x^{(n)} - \bar{x})^T v_1 \\ &= v_1^T \left(\frac{1}{N} \sum_{n=1}^N (x^{(n)} - \bar{x})(x^{(n)} - \bar{x})^T \right) v_1 = v_1^T \mathbf{S} v_1 \\ & \quad \text{s. t. } v_1^T v_1 = 1 \end{aligned}$$

First Principal Component

$$\max_v \frac{1}{N} \sum_{n=1}^N (v_1^T x^{(n)} - v_1^T \bar{x})^2 = v_1^T S v_1$$

s. t. $v_1^T v_1 = 1$

$$L(v_1, \lambda_1) = v_1^T S v_1 + \lambda_1 (1 - v_1^T v_1)$$

$$\frac{\partial L}{\partial v_1} = 0 \Rightarrow 2Sv_1 - 2\lambda_1 v_1 = 0$$
$$\Rightarrow Sv_1 = \lambda_1 v_1$$

Eigenvector with maximum eigenvalue maximizes the objective

Second Principal Component

$$\begin{aligned} \max_{v_2} \quad & v_2^T S v_2 \\ \text{s. t.} \quad & v_2^T v_2 = 1 \\ & v_2^T v_1 = 0 \end{aligned}$$

$$L(v_2, \lambda_2, \alpha) = v_2^T S v_2 + \lambda_2(1 - v_2^T v_2) - \alpha v_2^T v_1$$

Finding α :

$$\begin{aligned} \frac{\partial L}{\partial v_2} = 0 &\Rightarrow 2Sv_2 - 2\lambda_2 v_2 - \alpha v_1 = 0 \\ &\Rightarrow 2v_1^T S v_2 - 2\lambda_2 v_1^T v_2 - \alpha v_1^T v_1 = 0 \\ &\Rightarrow 2\lambda_1 v_1^T v_2 - 2\lambda_2 \times 0 - \alpha = 0 \\ &\Rightarrow \alpha = 0 \end{aligned}$$

Second Principal Component

$$\begin{aligned} \max_{v_2} & v_2^T S v_2 \\ \text{s. t.} & v_2^T v_2 = 1 \\ & v_2^T v_1 = 0 \end{aligned}$$

$$L(v_2, \lambda_2, \alpha) = v_2^T S v_2 + \lambda_2(1 - v_2^T v_2) - \alpha v_2^T v_1$$

Finding λ_2 :

$$\begin{aligned} \frac{\partial L}{\partial v_2} = 0 & \Rightarrow 2Sv_2 - 2\lambda_2 v_2 = 0 \\ & \Rightarrow Sv_2 = \lambda_2 v_2 \end{aligned}$$

λ_2 is the second largest eigenvalue

PCA Steps

- ▶ Input: $N \times d$ data matrix X (each row contain a d dimensional data point)
 - ▶ $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$
 - ▶ $\tilde{X} \leftarrow$ Mean value of data points is subtracted from rows of X
 - ▶ $S = \frac{1}{N} \tilde{X}^T \tilde{X}$ (Covariance matrix)
 - ▶ Calculate eigenvalue and eigenvectors of S
 - ▶ Pick d' eigenvectors corresponding to the largest eigenvalues and put them in the columns of $A = [\mathbf{v}_1, \dots, \mathbf{v}_{d'}]$
 - ▶ $X' = XA$
 - First PC
 - d'-th PC

Reconstruction

- Original data can be reconstructed using the PCs and the transformed data

$$\mathbf{x}' = \begin{bmatrix} \mathbf{v}_1^T \mathbf{x} \\ \vdots \\ \mathbf{v}_{d'}^T \mathbf{x} \end{bmatrix}$$

$$\mathbf{A} = [\mathbf{v}_1, \dots, \mathbf{v}_{d'}]$$

$$\mathbf{x}' = \mathbf{A}^T (\mathbf{x} - \bar{\mathbf{x}})$$

$$\Rightarrow \hat{\mathbf{x}} = \bar{\mathbf{x}} + \mathbf{A}\mathbf{x}' = \bar{\mathbf{x}} + \mathbf{A}\mathbf{A}^T (\mathbf{x} - \bar{\mathbf{x}})$$

if we use all PCs, then original data can be found without any error.

PCA on faces: Eigenfaces

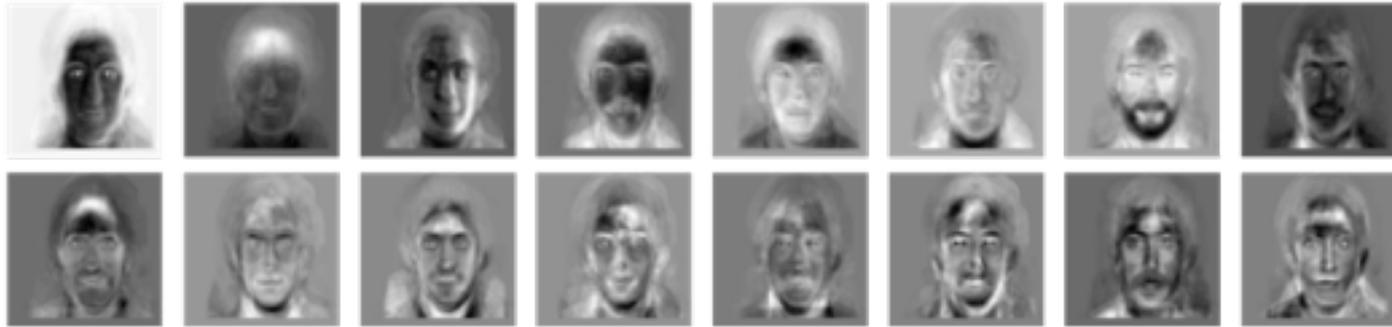


Display regular faces of dataset

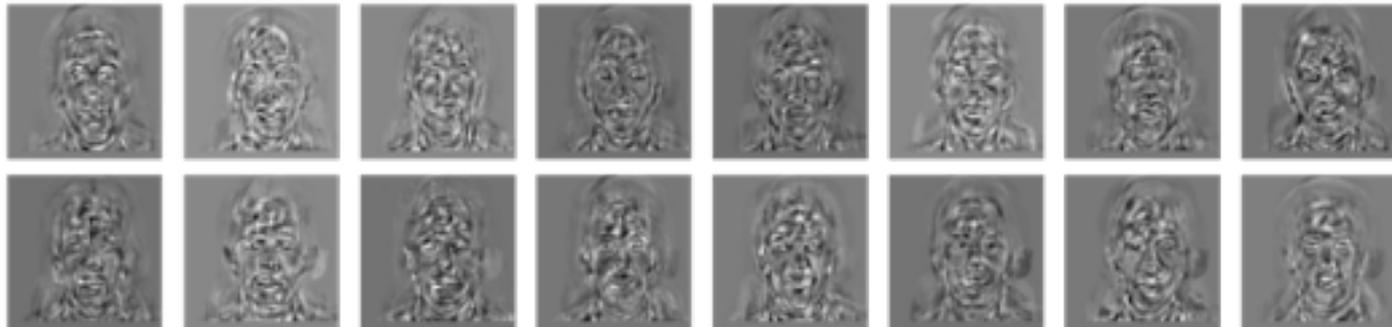


Mean face

Eigenfaces: PCs of faces



First few eigen faces **vs.** Last few eigenfaces



Eigenfaces: Reconstructing images

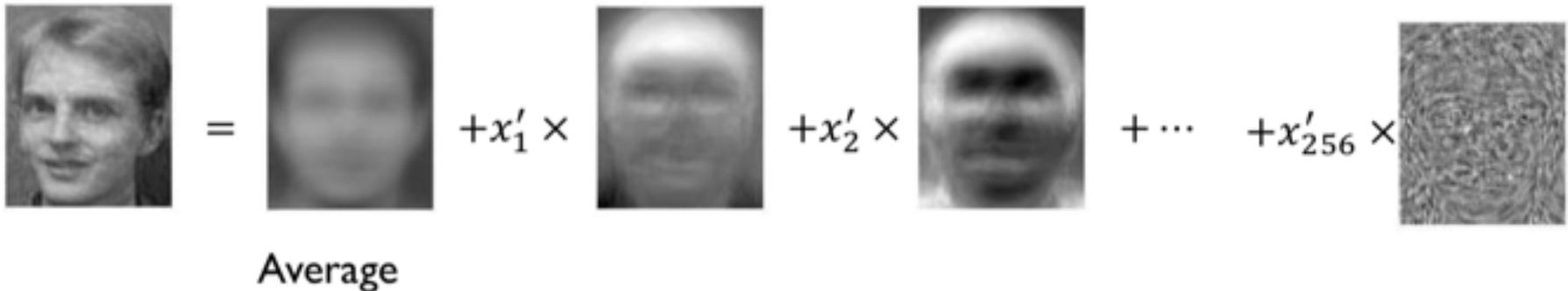


x is a $112 \times 92 = 10304$ dimensional vector containing intensity of the pixels of this image and $\tilde{x} = x - \bar{x}$

Feature vector = $[x'_1, x'_2, \dots, x'_{d'}]$

$x'_i = v_i^T \tilde{x} \rightarrow$ The projection of x on the i -th PC

$$\hat{x} = \bar{x} + \sum_{i=1}^{d'} x'_i \times v_i \qquad \hat{x} = \bar{x} + \sum_{i=1}^{d'} (v_i^T \tilde{x}) \times v_i$$



Eigenfaces: Reconstructing images

$d'=1$



$d'=2$



$d'=4$



$d'=8$



$d'=16$



$d'=32$



$d'=64$



$d'=128$



$d'=256$

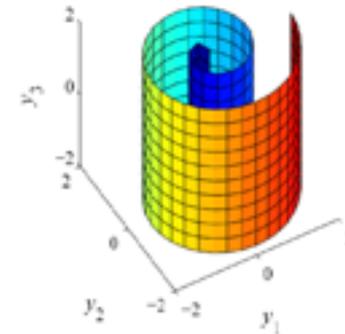
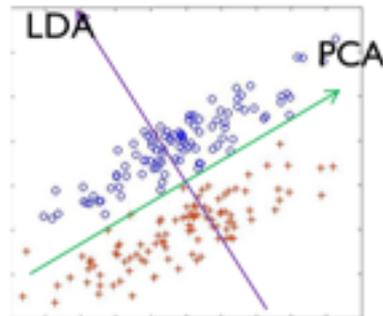
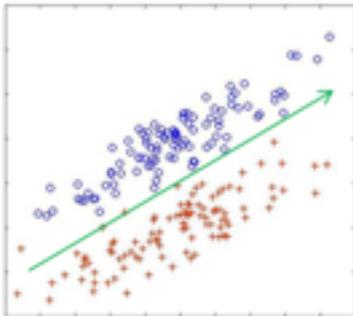


Original
Image



Pros & Cons

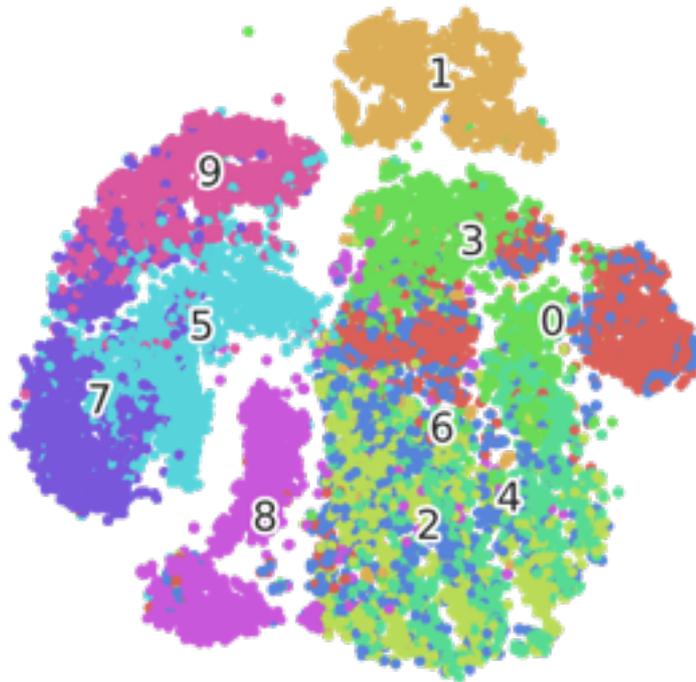
- No parameter tuning
- PCA is deterministic and fast
- Can be a preprocessing step specially for noise reduction in data
- However:
 - Data of different classes may not be separable after PCA
 - Distance among data and data topology may not be preserved



t-SNE

t-SNE

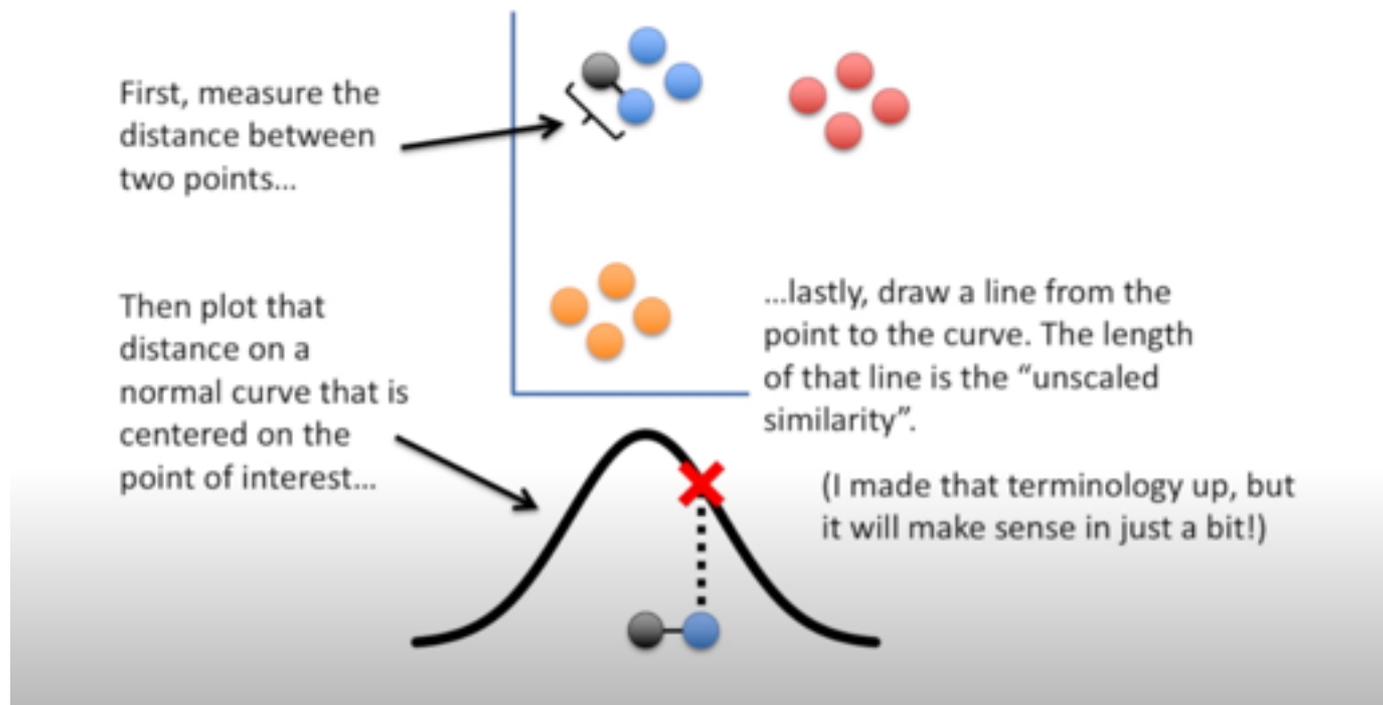
t-Distributed Stochastic Neighbor Embedding (t-SNE) is an **unsupervised**, **non-linear** technique primarily used for data exploration and **visualizing** high-dimensional data. The technique is a variation of Stochastic Neighbor Embedding (SNE).



Intuition

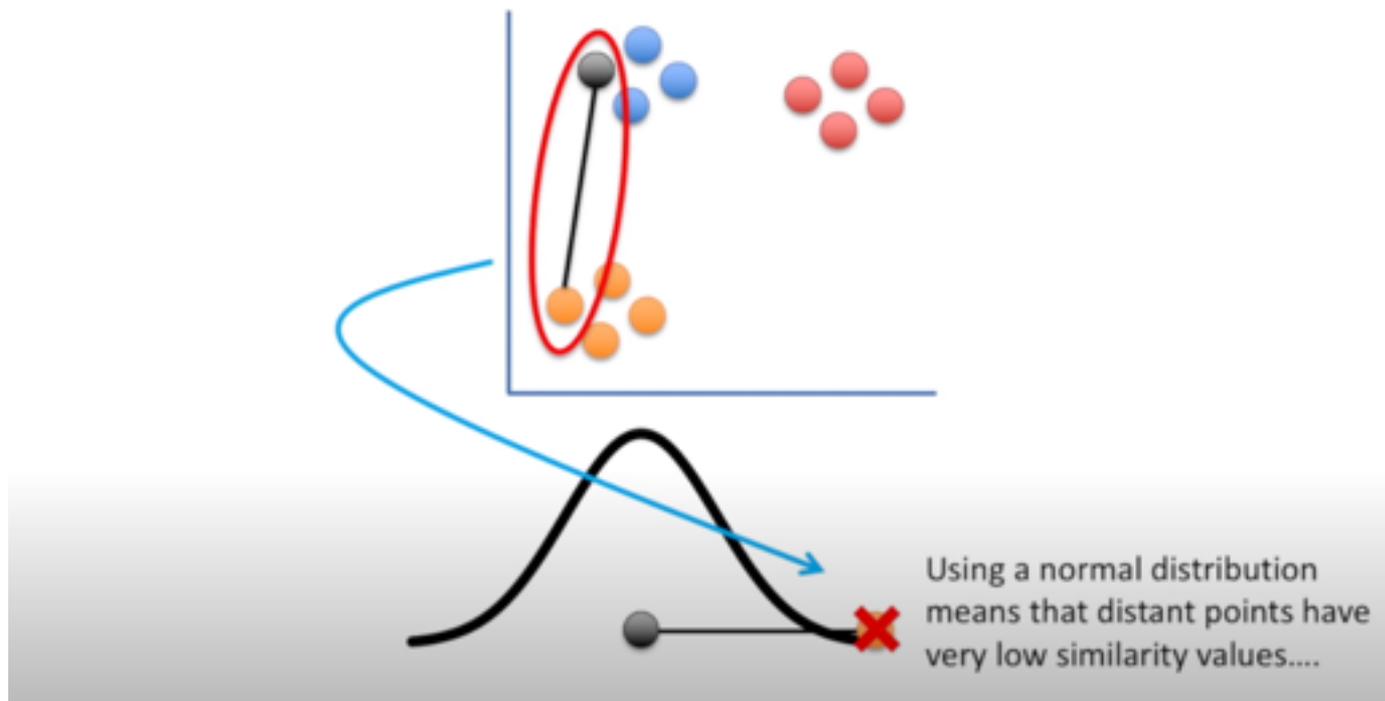
overview of working of t-SNE:

The algorithm starts by calculating the probability of similarity of points in high-dimensional space and calculating the probability of similarity of points in the corresponding low-dimensional space.



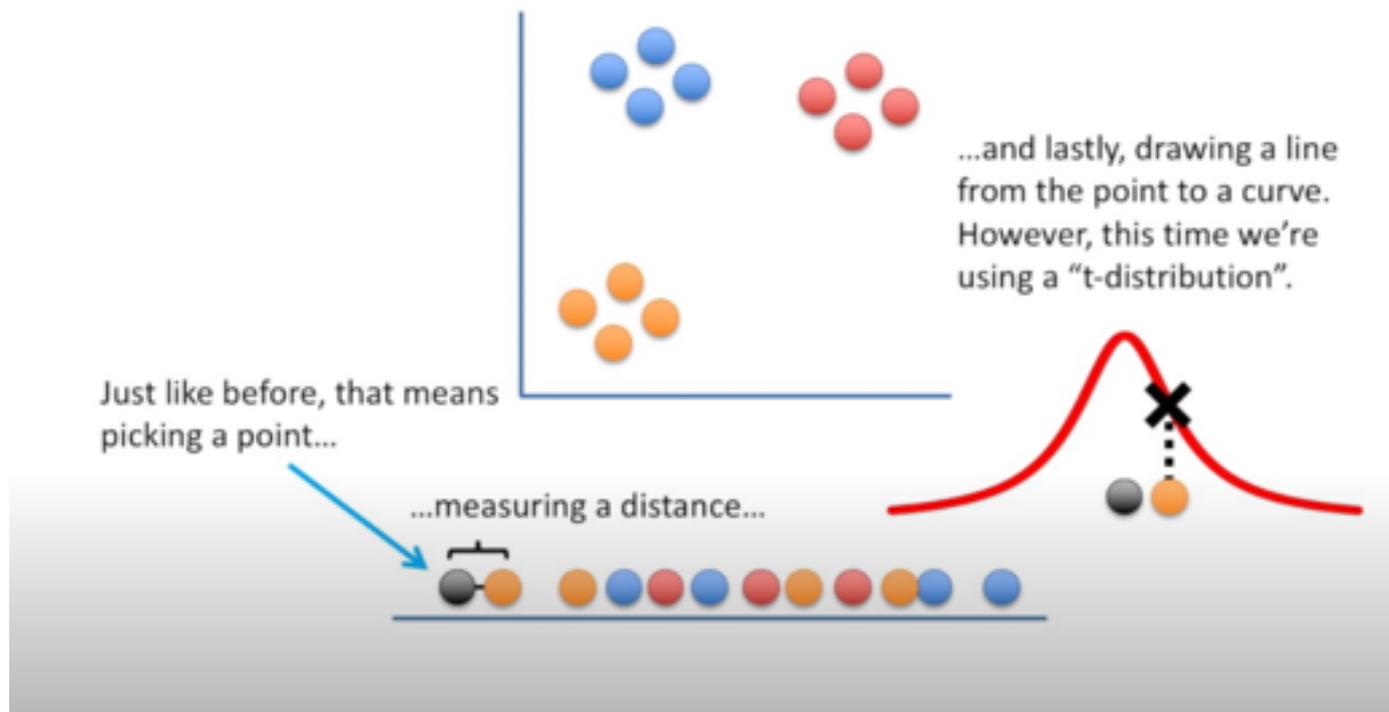
Intuition

The similarity of points is calculated as the conditional probability that a point A would choose point B as its neighbor if neighbors were picked in proportion to their probability density under a Gaussian (normal distribution) centered at A.

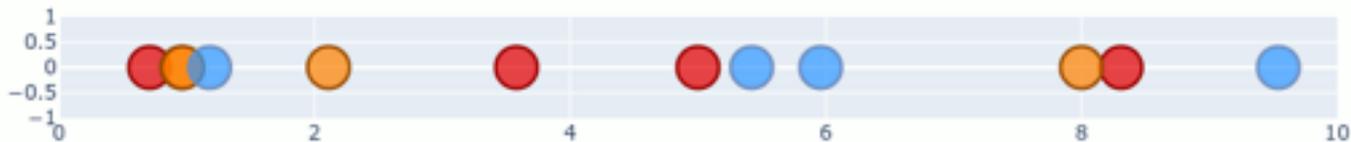
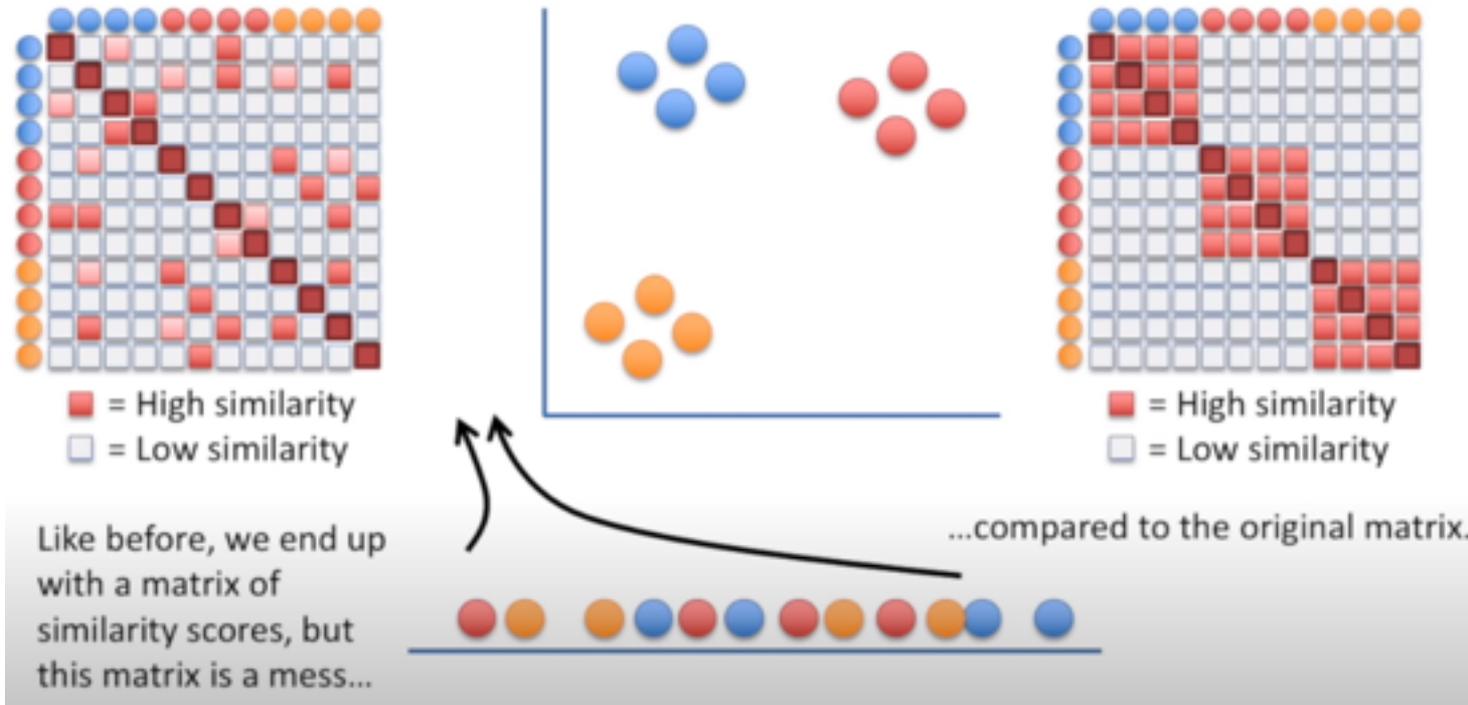


Intuition

It then tries to minimize the difference between these conditional probabilities (or similarities) in higher-dimensional and lower-dimensional space for a perfect representation of data points in lower-dimensional space.

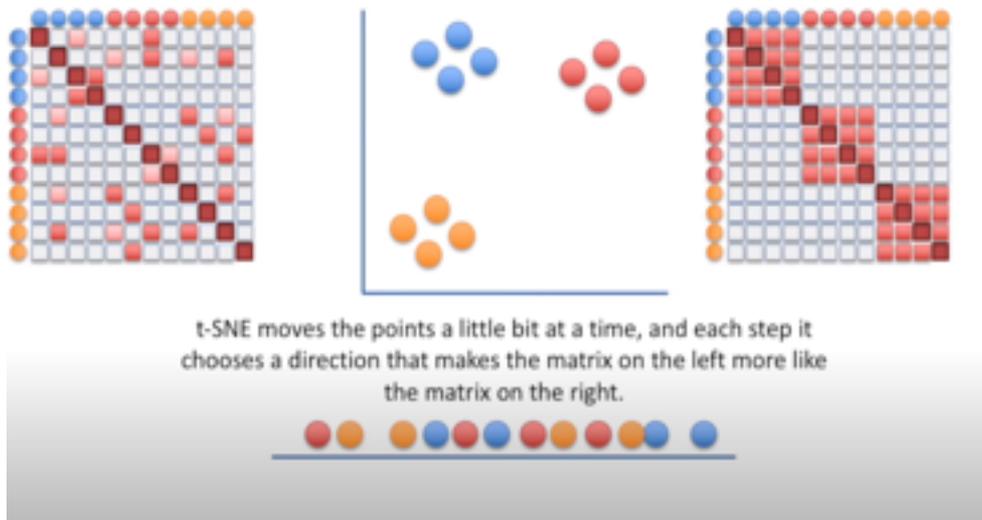


Intuition

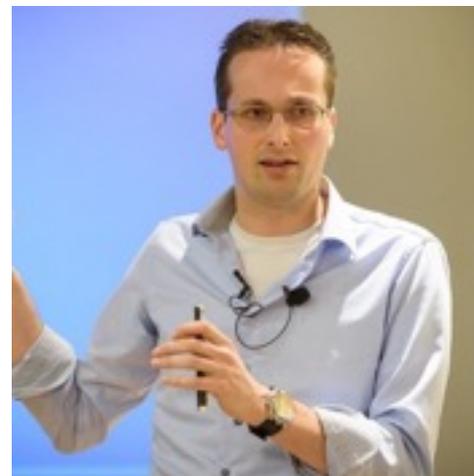
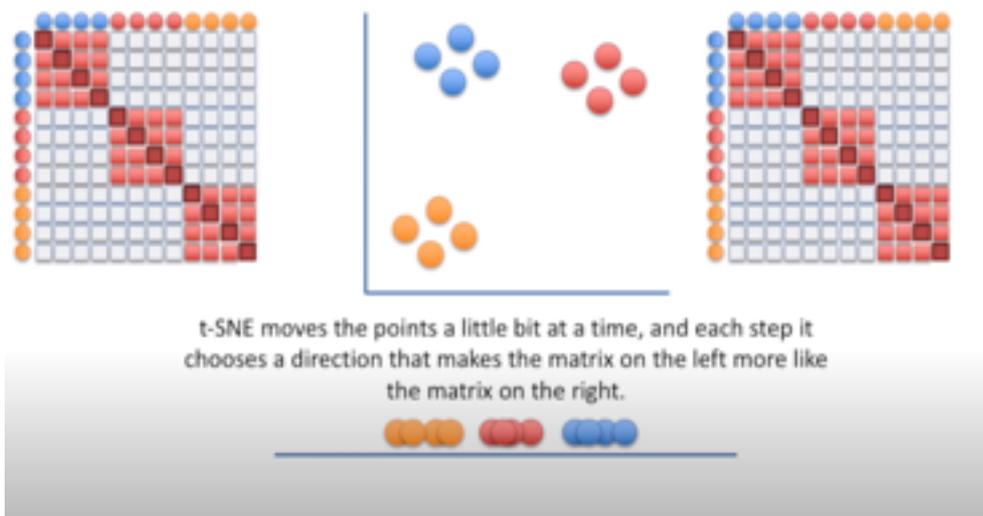


Intuition

In simpler terms, t-SNE gives you a feel or intuition of how the data is arranged in a high-dimensional space. It was developed by Laurens van der Maatens and Geoffrey Hinton in 2008.



Geoffrey Hinton



Laurens van der Maaten

SNE pair-wise similarities

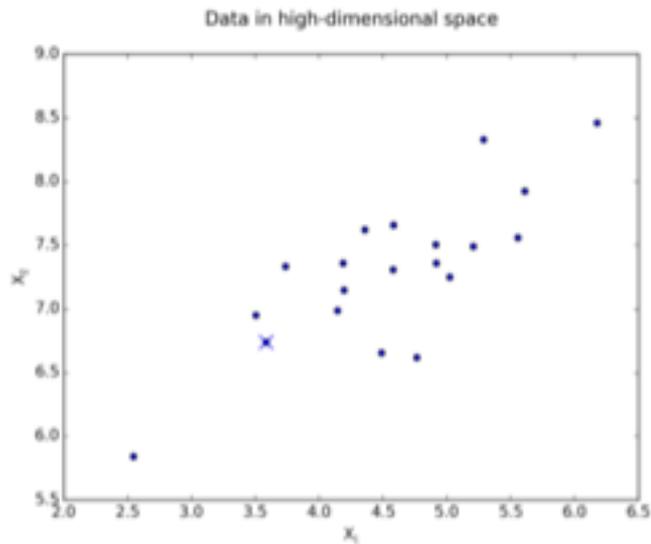
dimensionality reduction methods convert the high-dimensional data set $X = \{x_1, \dots, x_n\}$ into two or three-dimensional data $Y = \{y_1, \dots, y_n\}$ that can be displayed in a scatterplot. Stochastic Neighbor Embedding (SNE) starts by converting the high-dimensional Euclidean distances between datapoints into conditional probabilities that represent similarities.

$$P_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)} \quad p_{j|j} = 0$$

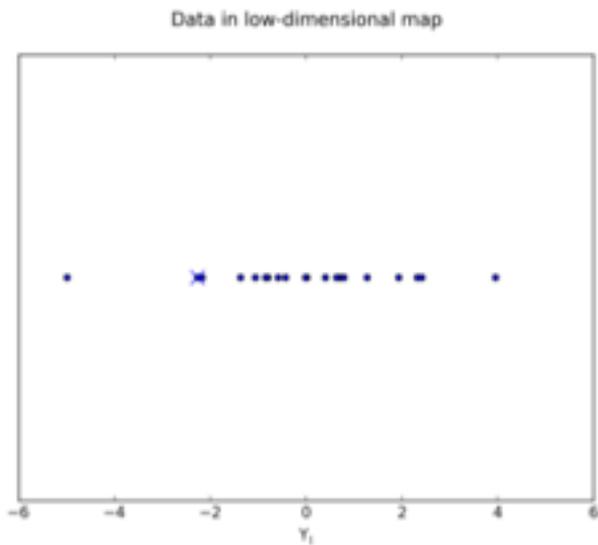
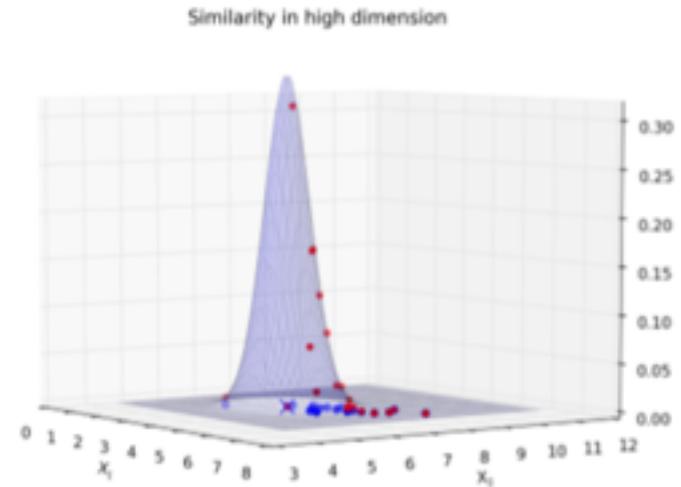
For nearby datapoints, $P_{j|i}$ is relatively high, whereas for widely separated datapoints, $P_{j|i}$ will be almost infinitesimal (for reasonable values of the variance of the Gaussian, σ). We set the variance of the Gaussian that is employed in the computation of the conditional probabilities $P_{j|i} = \frac{1}{\sqrt{2}}$.

$$q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)} \quad q_{j|j} = 0$$

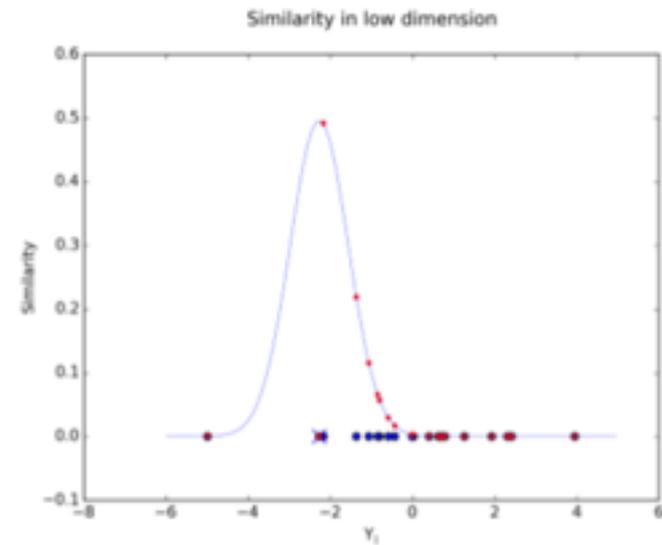
pair-wise similarities stay the same



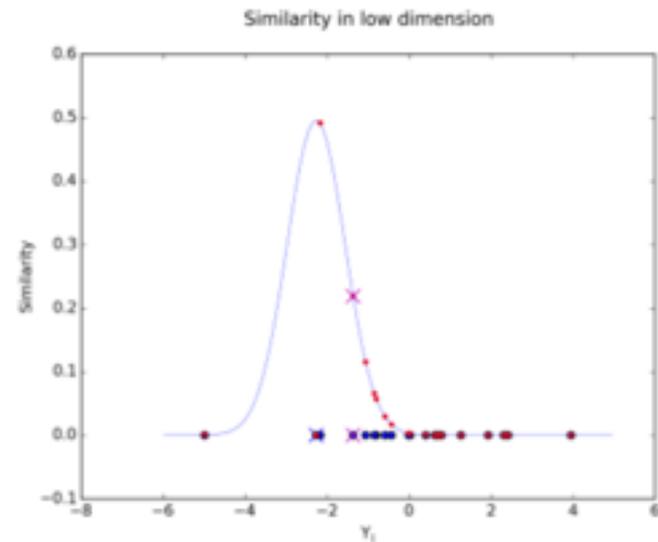
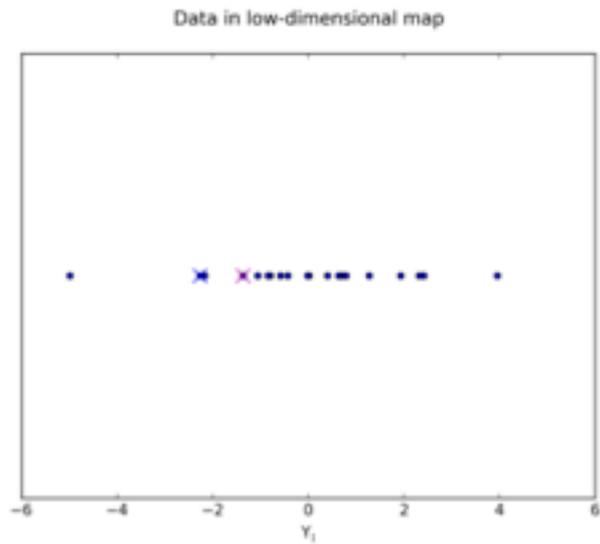
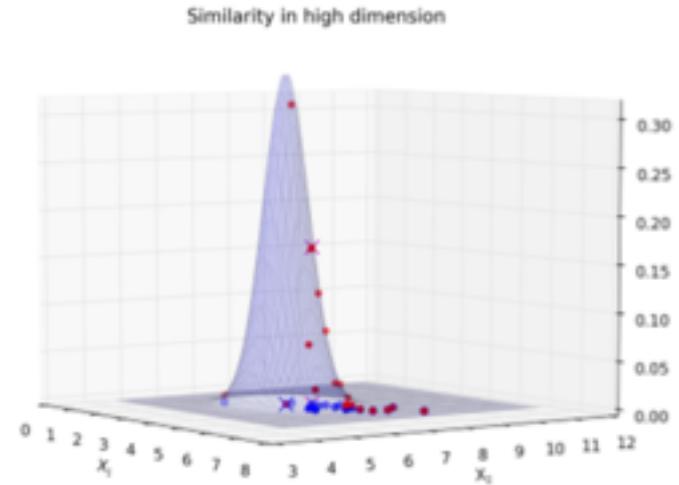
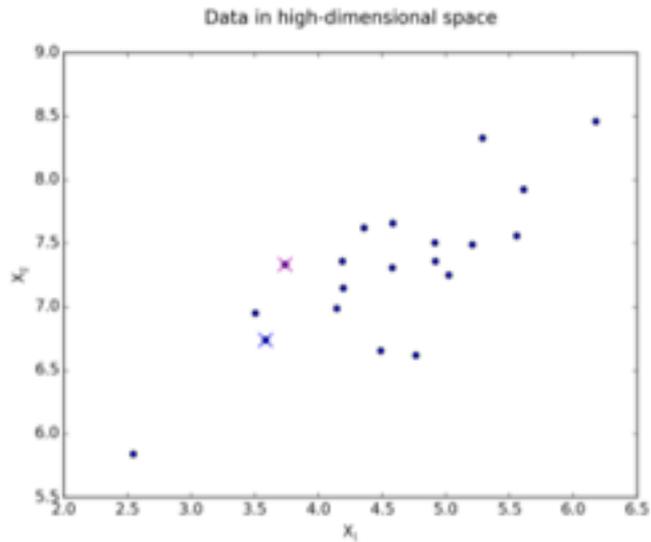
$$p_{j|i}$$
$$\Leftrightarrow$$



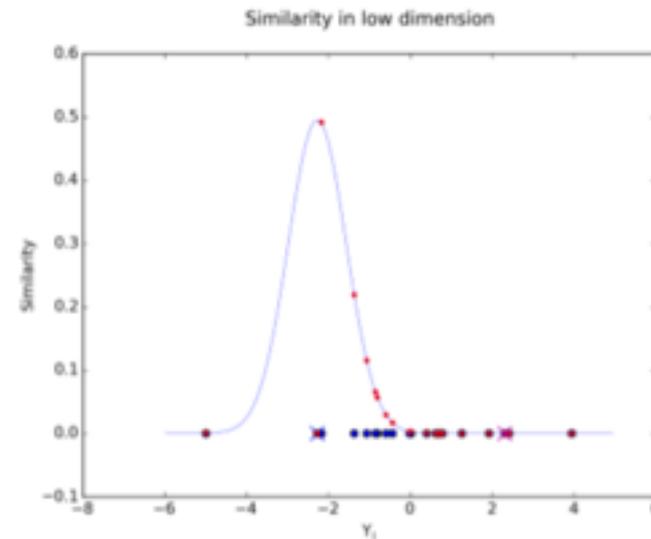
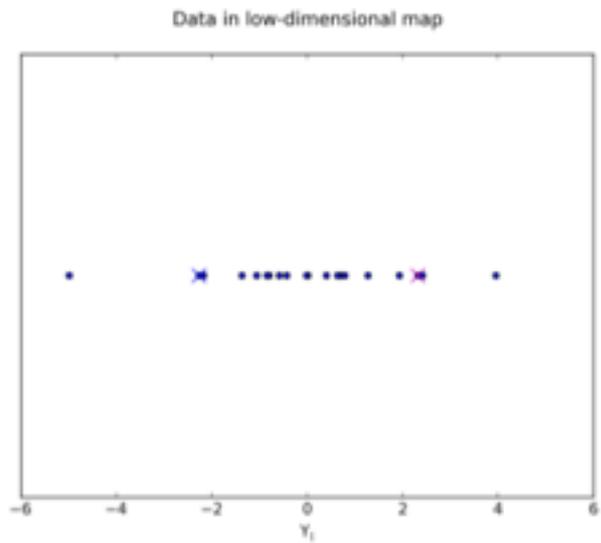
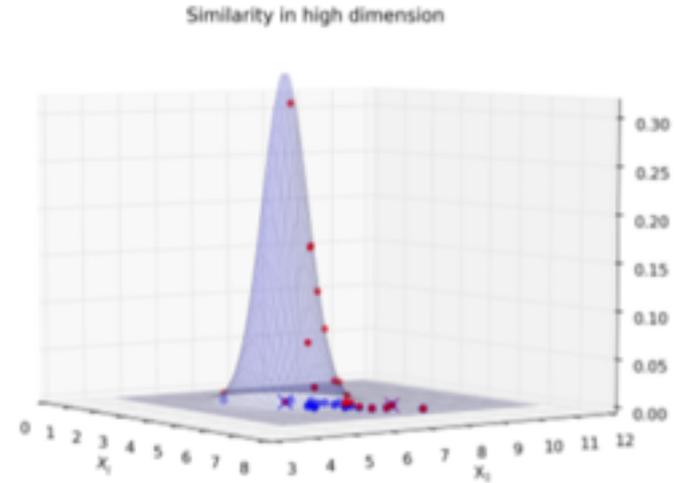
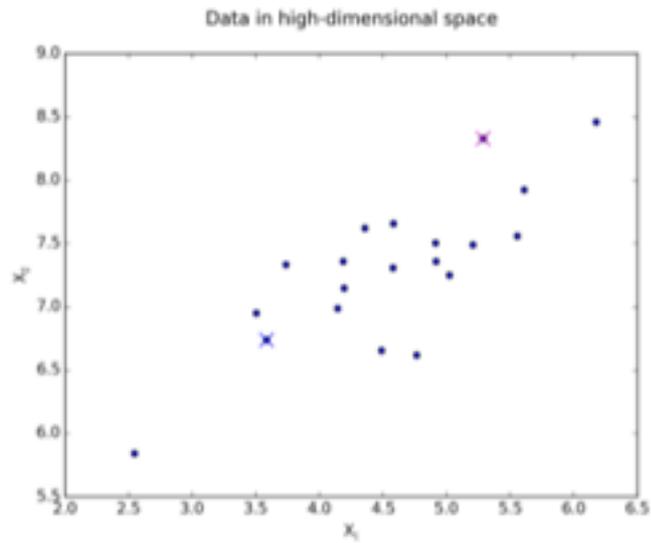
$$q_{j|i}$$
$$\Leftrightarrow$$



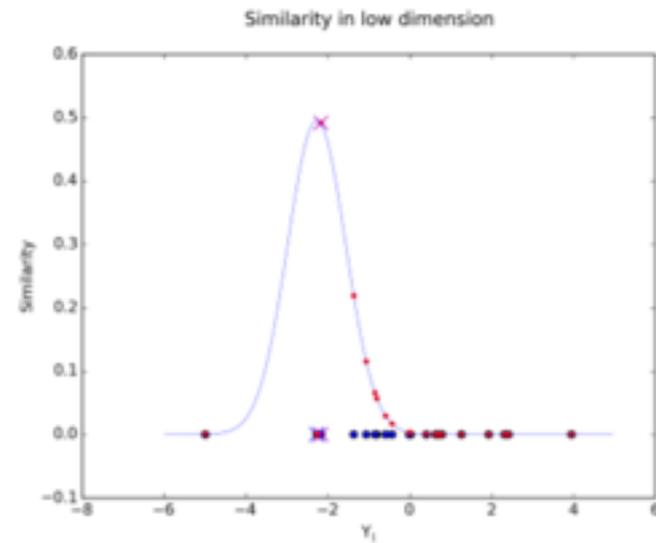
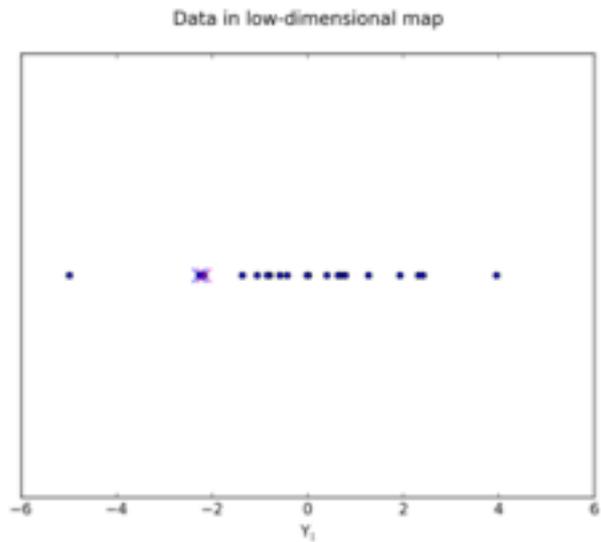
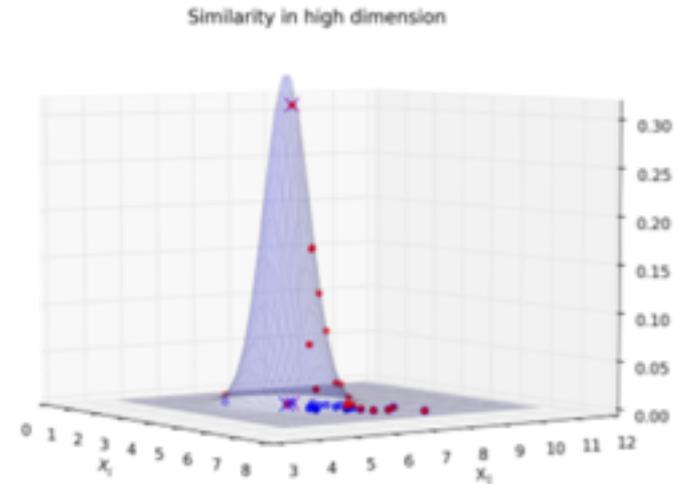
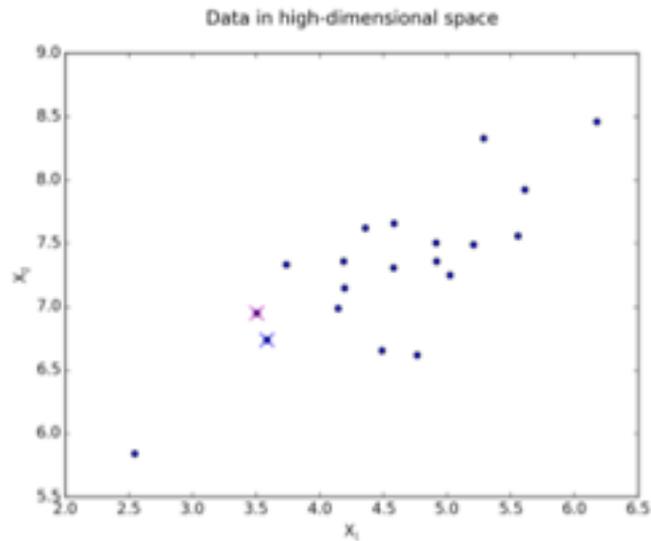
pair-wise similarities stay the same



pair-wise similarities stay the same



pair-wise similarities stay the same



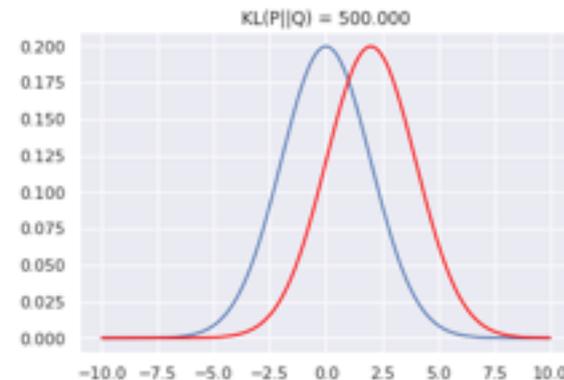
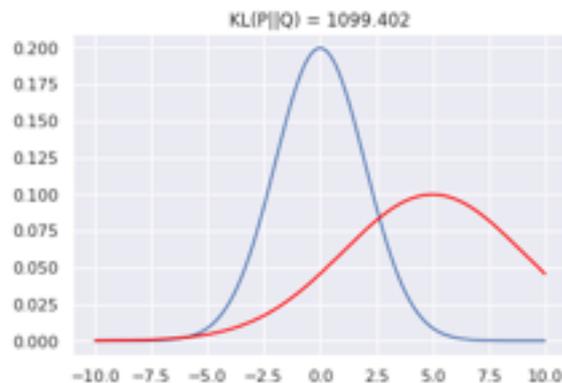
SNE

If the map points y_i and y_j correctly model the similarity between the high-dimensional datapoints x_i and x_j , the conditional probabilities $p_{j|i}$ and $q_{j|i}$ will be equal.

A natural measure of the faithfulness with which $q_{j|i}$ models $p_{j|i}$ is the **KullbackLeibler divergence** (which is in this case equal to the cross-entropy up to an additive constant). The Kullback-Leibler divergence, D_{KL} , is a measure of how one probability distribution is different from a second, reference probability distribution.

$$D_{KL}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

Lower the KL divergence value, the better we have matched the true distribution with our approximation.



SNE

SNE minimizes the sum of Kullback-Leibler divergences over all datapoints using a gradient descent method.

$$C = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$$

P_i represents the conditional probability distribution over all other datapoints given datapoint x_i , and Q_i represents the conditional probability distribution over all other map points given map point y_i .

It is not likely that there is a single value of σ_i that is optimal for all datapoints in the data set because the density of the data is likely to vary. In dense regions, a smaller value of σ_i is usually more appropriate than in sparser regions. SNE performs a binary search for the value of σ_i that produces a P_i with a fixed perplexity that is specified by the user.

$$H(P_i) = - \sum_j p_{j|i} \log_2 p_{j|i} \quad \text{Perp}(P_i) = 2^{H(P_i)}$$

Some questions

- **Why radial basis function (exponential)?**

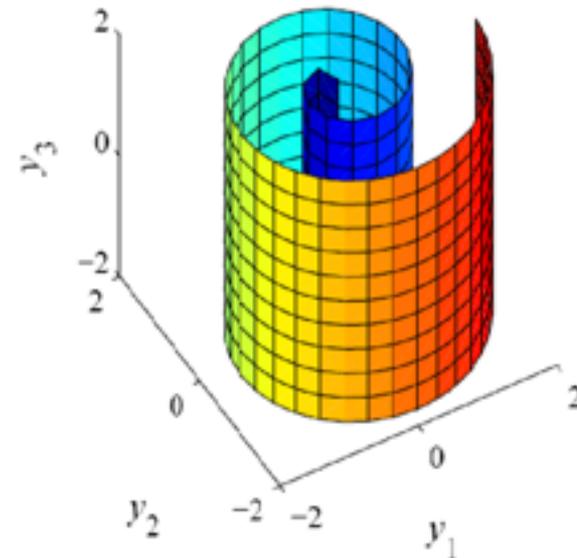
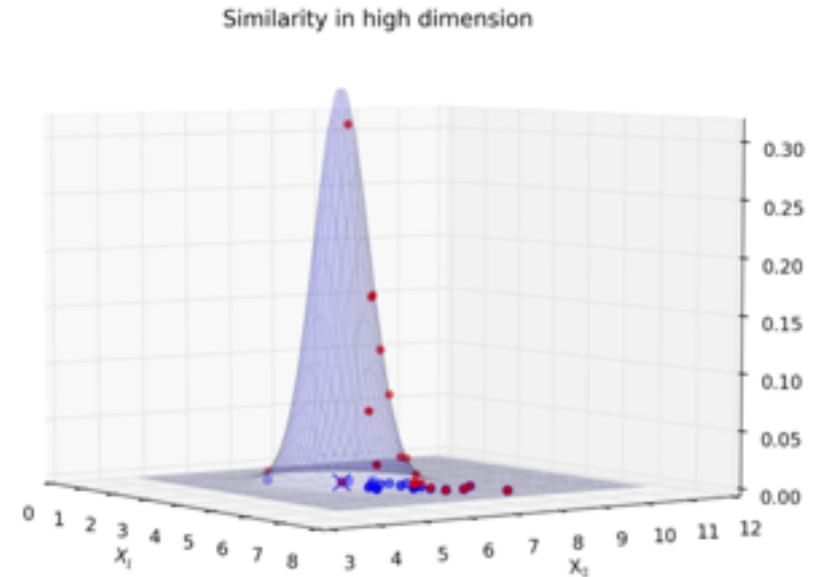
Focus on **local geometry**.

This is why t-SNE can be interpreted as **topology-based**.

- **Why probabilities?**

Small distance does not mean proximity on manifold.

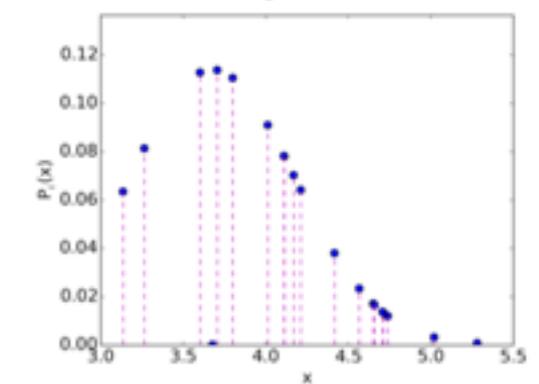
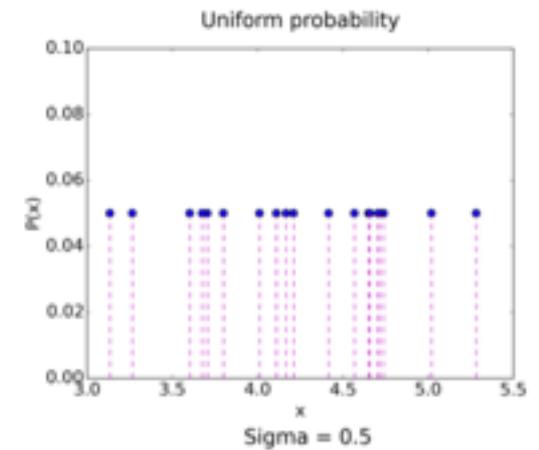
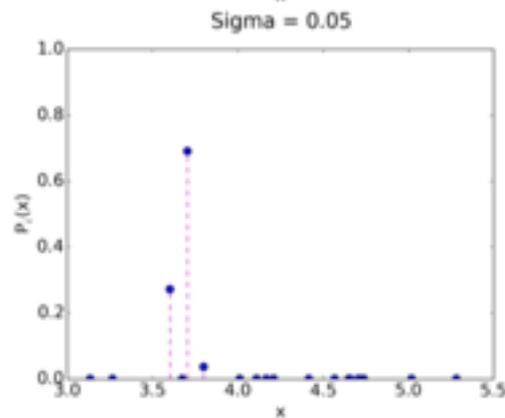
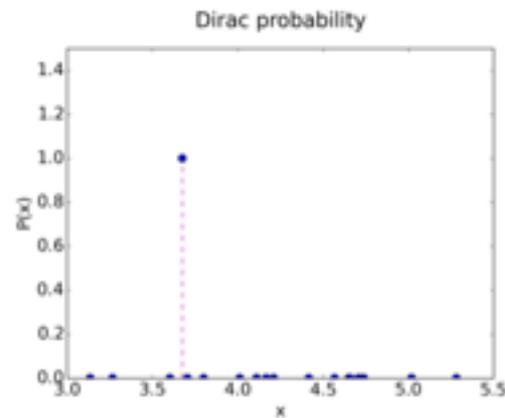
Probabilities are appropriate to model this **uncertainty**.



Some questions

- How do you choose σ ?

The entropy of P_i increases with σ .



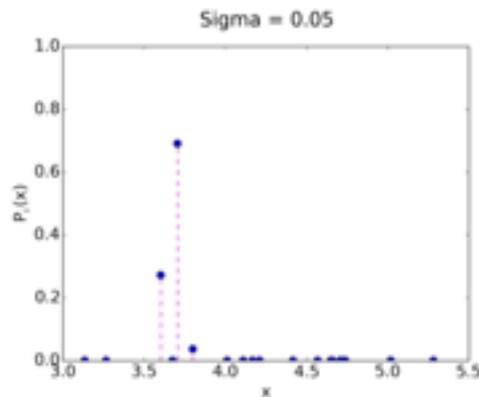
Some questions

- How do you choose σ ?

Perplexity, a smooth measure of the number of neighbors.

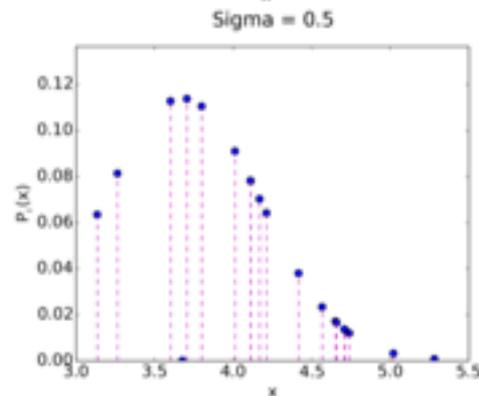
$$H(P) = -\sum P_i * \log_2 P_i$$

$$Prep(p) = 2^{H(p)}$$



\Rightarrow

Entropy of 1.055
Perplexity of 2.078



\Rightarrow

Entropy of 3.800
Perplexity of 13.929

SNE

The minimization of the cost function in Equation 2 is performed using a gradient descent method.

$$\frac{\delta C}{\delta y_i} = 2 \sum_j (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(y_i - y_j)$$

The current gradient is added to an exponentially decaying sum of previous gradients in order to determine the changes in the coordinates of the map points at each iteration of the gradient search.

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

t-SNE

In this section, we first discuss the symmetric version of SNE

As an alternative to minimizing the sum of the Kullback-Leibler divergences between the conditional probabilities $P_{j|i}$ and $Q_{j|i}$, it is also possible to minimize a single Kullback-Leibler divergence between a joint probability distribution, P , in the high-dimensional space and a joint probability distribution, Q , in the low-dimensional space

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

We refer to this type of SNE as symmetric SNE, because it has the property that $p_{ij} = p_{ji}$ and $q_{ij} = q_{ji}$ for $\forall i, j$.

$$q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq l} \exp(-\|y_k - y_l\|^2)} \quad p_{ij} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2 / 2\sigma^2)}$$

t-SNE

But this causes problems when a high-dimensional datapoint x_i is an outlier.

We set $P_{ij} = \frac{P_{ij} + P_{ji}}{2n}$. This ensures that $\sum_j P_{ij} > \frac{1}{2n}$ for all datapoints x_i , as a result of which each datapoint x_i makes a significant contribution to the cost function.

$$\frac{\delta C}{\delta y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j)$$

The main advantage of the symmetric version of SNE is the simpler form of its gradient, which is faster to compute.

SNE

In t-SNE, we employ a Student t-distribution with one degree of freedom as the heavy-tailed distribution in the low-dimensional map.

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

The gradient of the Kullback-Leibler divergence between P and the Student-t based joint probability distribution Q:

$$\frac{\delta C}{\delta y_i} = 4 \sum_j (p_{ij} - q_{ij})(1 + \|y_i - y_j\|^2)^{-1}(y_i - y_j)$$

- What are the differences between SNE and t-SNE?
- Why we should use t-SNE instead of SNE?

Pseudo Code

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)

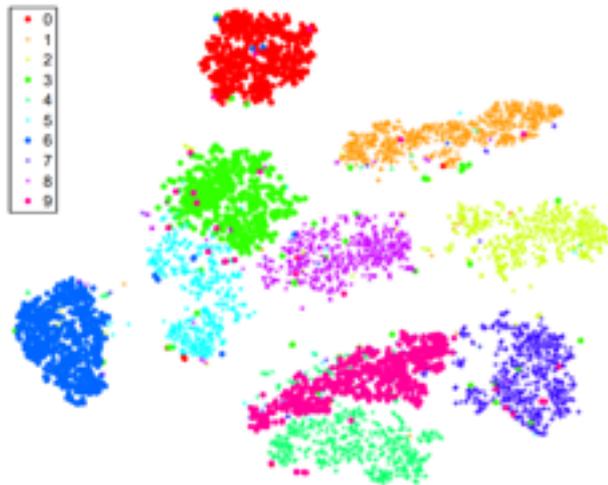
 compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)

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

end

end

Visualizations of 6,000 handwritten digits from the MNIST data set.



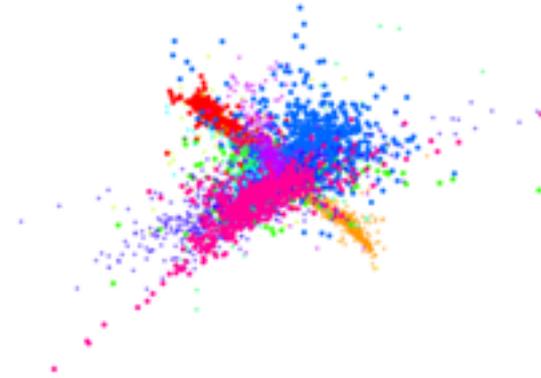
(a) Visualization by t-SNE.



(b) Visualization by Sammon mapping.

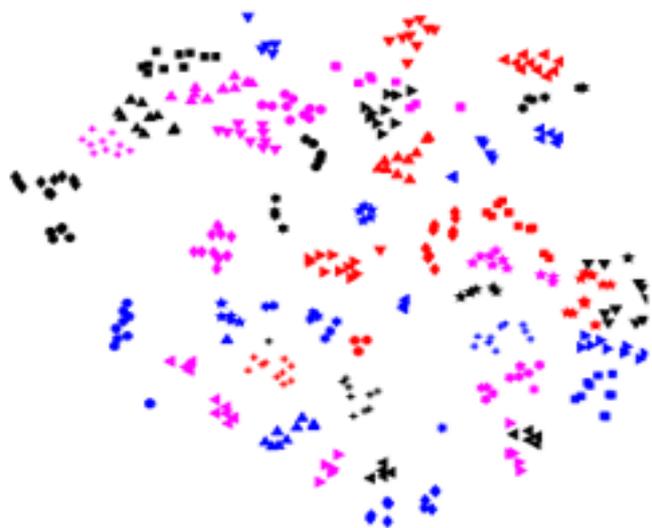


(a) Visualization by Isomap.



(b) Visualization by LLE.

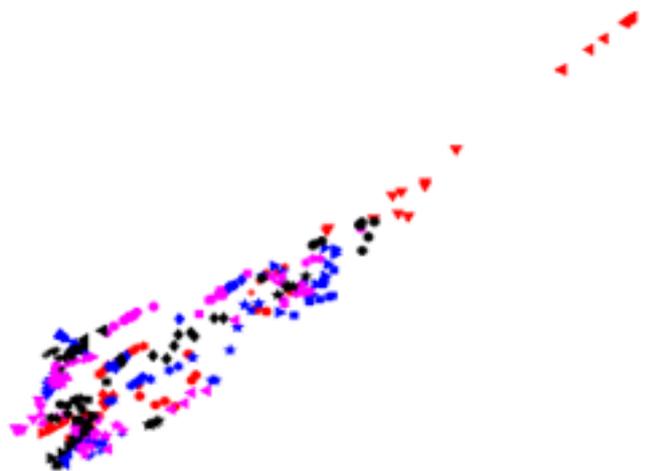
Visualizations of the Olivetti faces data set.



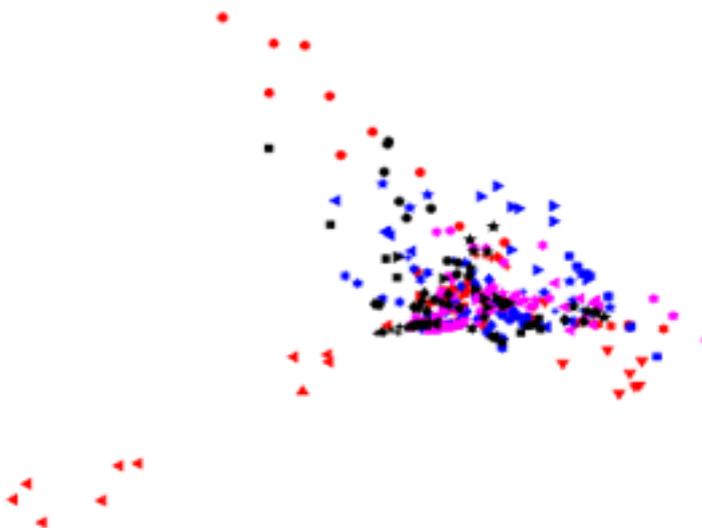
(a) Visualization by t-SNE.



(b) Visualization by Sammon mapping.



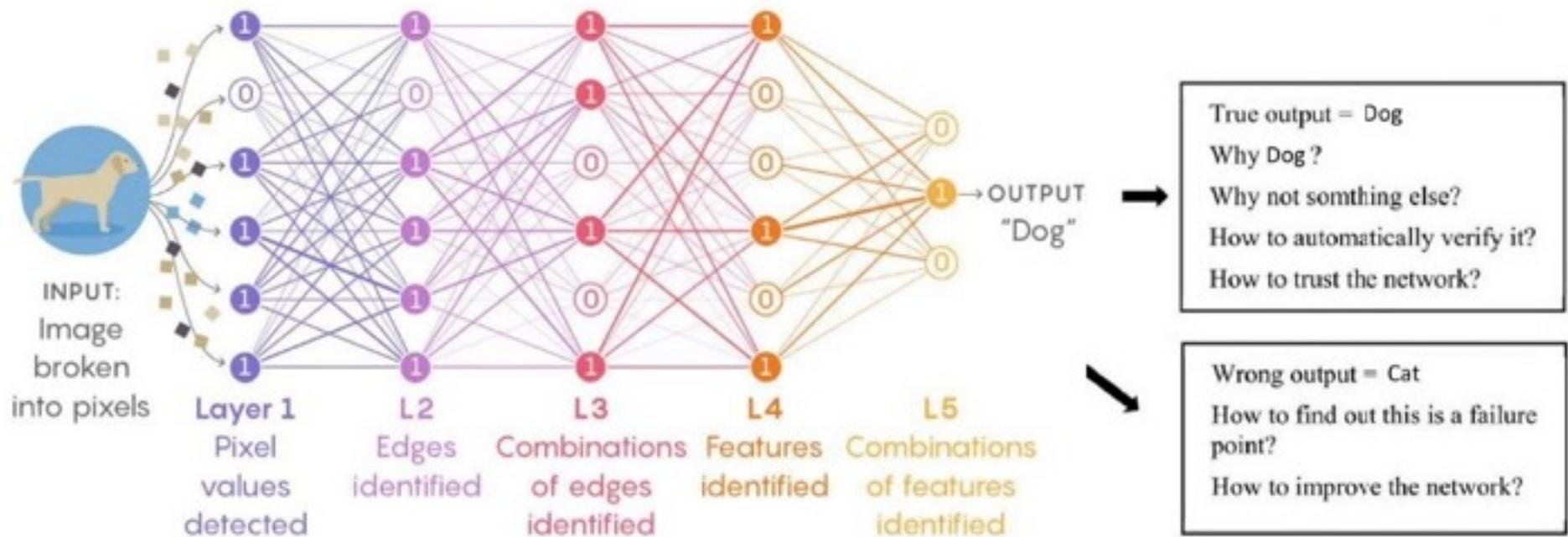
(c) Visualization by Isomap.



(d) Visualization by LLE.

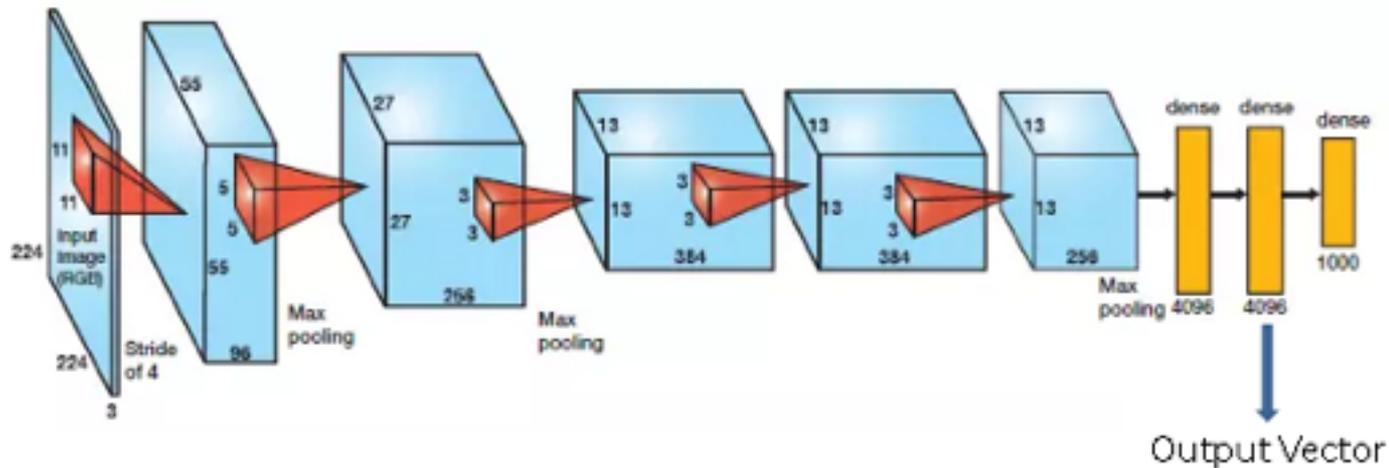
Important Application

The Deep Learning Algorithms mimic the cognitive capabilities of the human brain and the Interpretability of Deep Learning Algorithms has been an important area of research since the start of the Deep Learning era



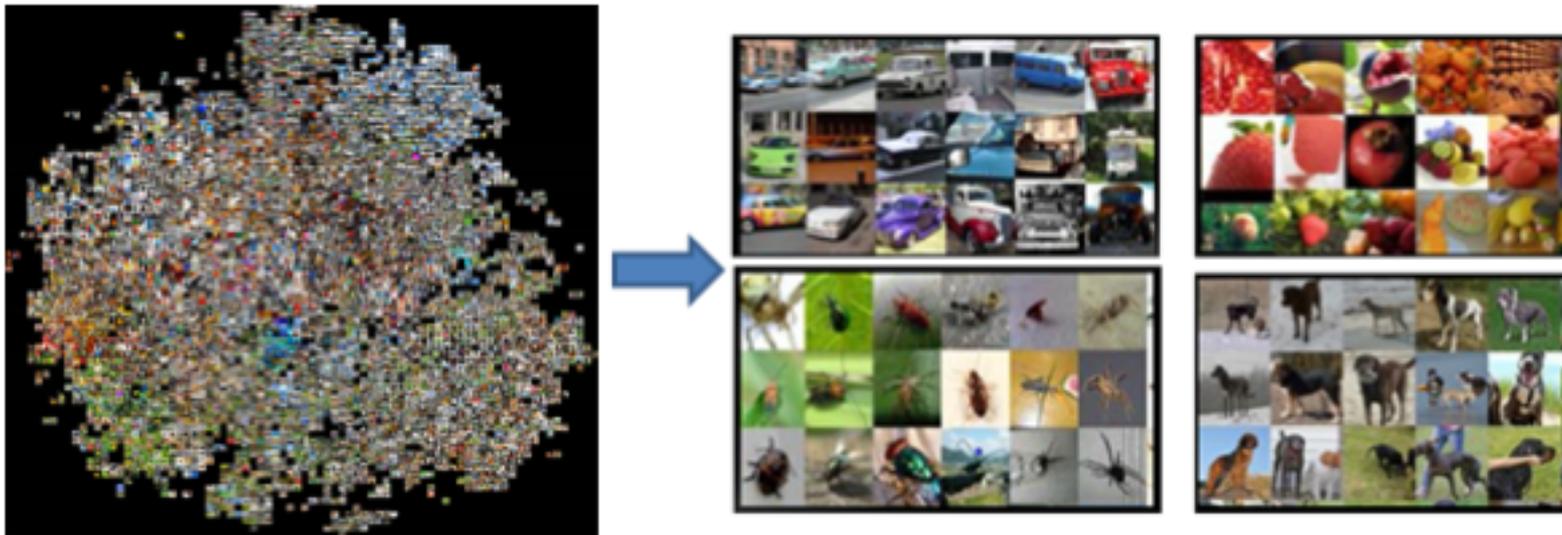
Important Application

ConvNets can be interpreted as gradually transforming the images into a representation in which the classes are separable by a linear classifier. We can get a rough idea about the topology of this space by embedding images into low dimensions and t-SNE is one powerful technique that can be used for embedding high-dimensional vectors in a low-dimensional space while preserving the pairwise distances of the points.



Important Application

To produce an embedding, take a set of images and Forward propagating each image through a trained ConvNet to extract a vector for each image class (e.g. in AlexNet the 4096-dimensional vector right before the classifier can be used for embedding) and then plug the vectors into t-SNE and get a 2-dimensional vector for each image.



Important Application

Images that are nearby each other are also close in the CNN representation space, which implies that CNN “sees” them as being very similar.



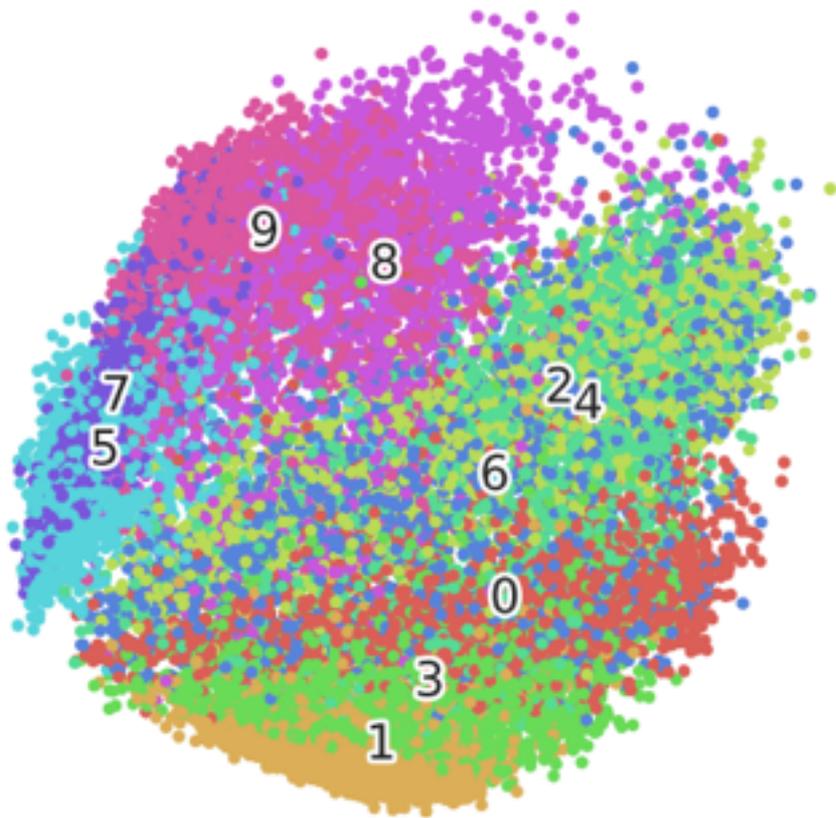
Cluster containing black cats.



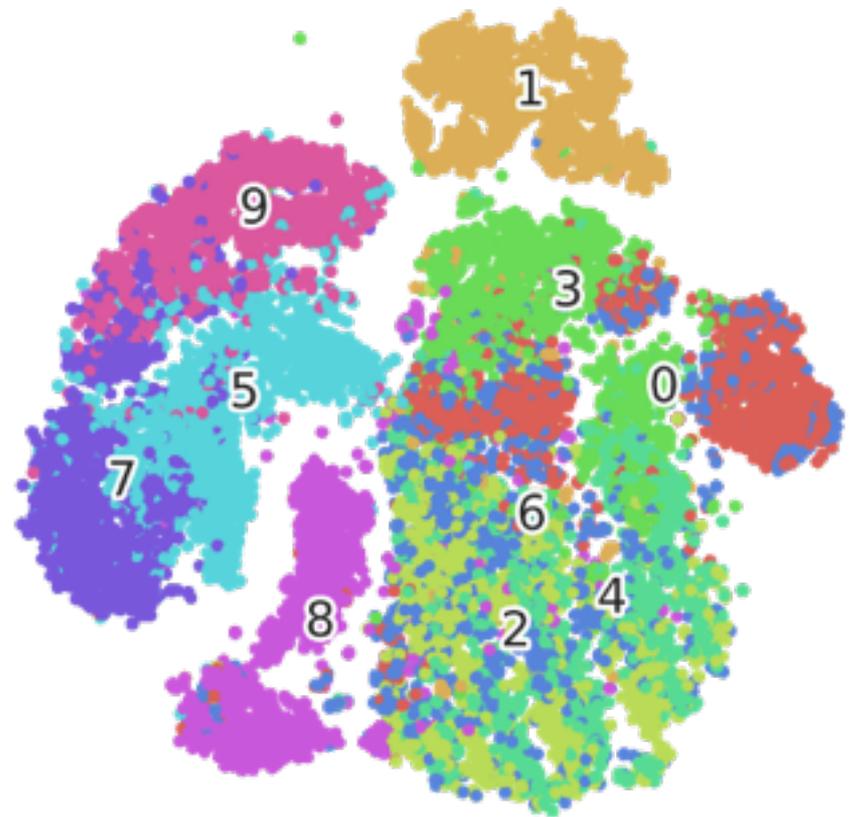
Cluster containing gray cats.

PCA vs t-SNE

you will see the output of PCA on the Fashion-MNIST dataset to compare it with t-SNE.



PCA



t-SNE

PCA vs t-SNE

some key differences between PCA and t-SNE can be noted as follows:

- t-SNE is computationally expensive and can take several hours on million-sample datasets where PCA will finish in seconds or minutes.
- PCA it is a mathematical technique, but t-SNE is a probabilistic one.
- Linear dimensionality reduction algorithms, like PCA, concentrate on placing dissimilar data points far apart in a lower dimension representation. But in order to represent high dimension data on low dimension, non-linear manifold, it is essential that similar data points must be represented close together, which is something t-SNE does not PCA.
- Sometimes in t-SNE different runs with the same hyperparameters may produce different results hence multiple plots must be observed before making any assessment with t-SNE, while this is not the case with PCA.
- Since PCA is a linear algorithm, it will not be able to interpret the complex polynomial relationship between features while t-SNE is made to capture exactly that.

References

- CE-717: Machine Learning _ Sharif University of Technology _ Fall 2020 _ Soleymani
- www.en.wikipedia.org/wiki/Principal_component_analysis
- www.cs.tau.ac.il/~rshamir/abdbm/pres/17/PCA.pdf
- <https://heartbeat.fritz.ai/understanding-the-mathematics-behind-principal-component-analysis-efd7c9ff0bb3>
- <https://www.youtube.com/watch?v=NEaUSP4YerM>
- https://en.wikipedia.org/wiki/T-distributed_stochastic_neighbor_embedding
- <https://towardsdatascience.com/t-sne-clearly-explained-d84c537f53a>
- <https://www.datacamp.com/community/tutorials/introduction-t-sne>
- <https://www.researchgate.net/post/Does-Random-Projection-have-advantages-over-PCA>
- <https://medium.com/@lwj.liuwenjing/how-to-get-eigenfaces-a9caeeba8767>
- <https://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf?fbclid=IwAR0Bgg1eA5TFmqOZeCQXsloL6PKrVXUFaskUKtg6yBhVXAFFvZA6yQiYx-M>
- <https://medium.com/analytics-vidhya/deep-learning-visualization-and-interpretation-of-neural-networks-2f3f82f501c5>

Any Questions?!