

Machine Learning 2

DS 4420 - Spring 2020

Dimensionality reduction 2

Byron C Wallace



Today

- A bit of wrap up on PCA
- **Then:** Non-linear dimensionality reduction! (SNE/t-SNE)

In Sum: Principal Component Analysis

Data

$$\mathbf{X} = \begin{pmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ | & & | \end{pmatrix} \in \mathbb{R}^{d \times n}$$

Eigenvectors of Covariance

$$\mathbf{C} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top = \frac{1}{n} \mathbf{X} \mathbf{X}^\top$$

$$\mathbf{C} \mathbf{u}_j = \lambda_j \mathbf{u}_j$$

$$\mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \cdots & \\ & & & \lambda_d \end{pmatrix}$$

Idea: Take **top- k** eigenvectors to maximize variance

Why?

$$\mathbf{C} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top = \frac{1}{n} \mathbf{X} \mathbf{X}^\top$$

$$\mathbf{C} \mathbf{u}_j = \lambda_j \mathbf{u}_j$$

Idea: Take **top- k** eigenvectors to maximize variance

Last time, we saw that we can derive this by *maximizing the variance in the compressed space*

Why?

$$\mathbf{C} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top = \frac{1}{n} \mathbf{X} \mathbf{X}^\top$$

$$\mathbf{C} \mathbf{u}_j = \lambda_j \mathbf{u}_j$$

Idea: Take **top- k** eigenvectors to maximize variance

Last time, we saw that we can derive this by *maximizing the variance in the compressed space*

Can also motivate by explicitly minimizing *reconstruction error*

Minimizing *reconstruction error*

Getting the eigenvalues, two ways

- Direct eigenvalue decomposition of the covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top$$

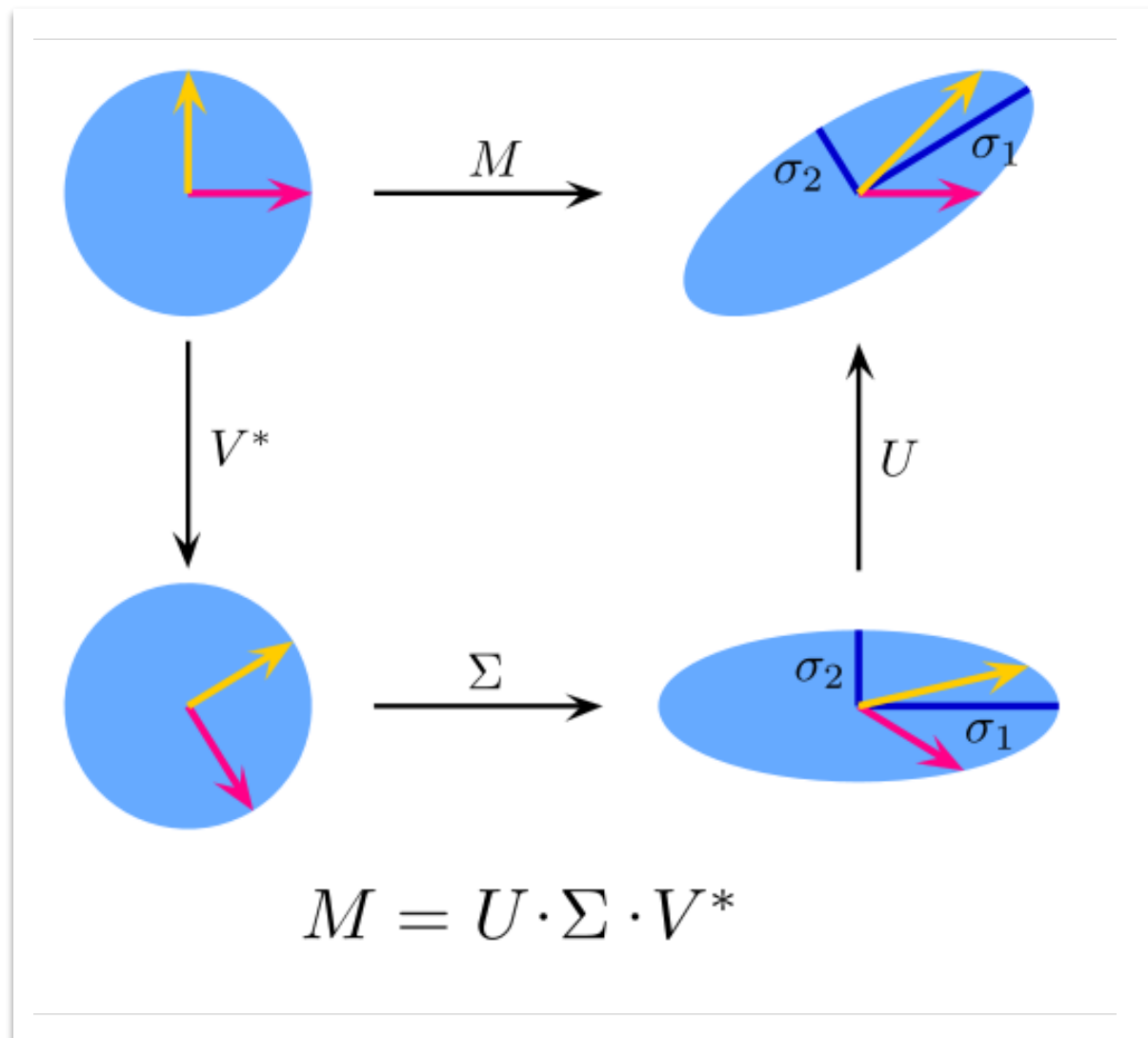
Getting the eigenvalues, two ways

- Direct eigenvalue decomposition of the covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top$$

- Singular Value Decomposition (SVD)

Singular Value Decomposition



Idea: Decompose the $d \times n$ matrix \mathbf{X} into

1. A $n \times n$ basis \mathbf{V} (unitary matrix)
2. A $d \times n$ matrix $\mathbf{\Sigma}$ (diagonal projection)
3. A $d \times d$ basis \mathbf{U} (unitary matrix)

$$\mathbf{X} = \mathbf{U}_{d \times d} \mathbf{\Sigma}_{d \times n} \mathbf{V}_{n \times n}^{\top}$$

1. Rotation

$$V^T \vec{x} = \sum_{i=1}^n \langle \vec{v}_i, \vec{x} \rangle \vec{e}_i$$

2. Scaling

$$SV^T \vec{x} = \sum_{i=1}^n s_i \langle \vec{v}_i, \vec{x} \rangle \vec{e}_i$$

3. Rotation

$$USV^T \vec{x} = \sum_{i=1}^n s_i \langle \vec{v}_i, \vec{x} \rangle \vec{u}_i$$

SVD for PCA

$$\underbrace{\mathbf{X}}_{D \times N} = \underbrace{\mathbf{U}}_{D \times D} \underbrace{\mathbf{\Sigma}}_{D \times N} \underbrace{\mathbf{V}^\top}_{N \times N}$$

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \underbrace{\mathbf{V}^\top \mathbf{V}}_{=\mathbf{I}_N} \mathbf{\Sigma}^\top \mathbf{U}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^\top \mathbf{U}^\top$$

SVD for PCA

$$\underbrace{\mathbf{X}}_{D \times N} = \underbrace{\mathbf{U}}_{D \times D} \underbrace{\mathbf{\Sigma}}_{D \times N} \underbrace{\mathbf{V}^\top}_{N \times N}$$

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \underbrace{\mathbf{V}^\top \mathbf{V}}_{=\mathbf{I}_N} \mathbf{\Sigma}^\top \mathbf{U}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^\top \mathbf{U}^\top$$

It turns out the columns of \mathbf{U} are the eigenvectors of $\mathbf{X} \mathbf{X}^\top$

Computing PCA

Method 1: eigendecomposition

\mathbf{U} are eigenvectors of covariance matrix $C = \frac{1}{n}\mathbf{X}\mathbf{X}^\top$

Computing C already takes $O(nd^2)$ time (very expensive)

Method 2: singular value decomposition (SVD)

Find $\mathbf{X} = \mathbf{U}_{d \times d} \Sigma_{d \times n} \mathbf{V}_{n \times n}^\top$

where $\mathbf{U}^\top \mathbf{U} = I_{d \times d}$, $\mathbf{V}^\top \mathbf{V} = I_{n \times n}$, Σ is diagonal

Computing top k singular vectors takes only $O(ndk)$

Computing PCA

Method 1: eigendecomposition

\mathbf{U} are eigenvectors of covariance matrix $C = \frac{1}{n}\mathbf{X}\mathbf{X}^\top$

Computing C already takes $O(nd^2)$ time (very expensive)

Method 2: singular value decomposition (SVD)

Find $\mathbf{X} = \mathbf{U}_{d \times d} \Sigma_{d \times n} \mathbf{V}_{n \times n}^\top$

where $\mathbf{U}^\top \mathbf{U} = I_{d \times d}$, $\mathbf{V}^\top \mathbf{V} = I_{n \times n}$, Σ is diagonal

Computing top k singular vectors takes only $O(ndk)$

Relationship between eigendecomposition and SVD:

Left singular vectors are principal components ($C = \mathbf{U}\Sigma^2\mathbf{U}^\top$)

Eigen-faces [Turk & Pentland 1991]

- d = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a face image
- \mathbf{x}_{ji} = intensity of the j -th pixel in image i

Eigen-faces [Turk & Pentland 1991]

- d = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a face image
- \mathbf{x}_{ji} = intensity of the j -th pixel in image i

$$\mathbf{X}_{d \times n} \approx \mathbf{U}_{d \times k} \mathbf{Z}_{k \times n}$$

$\left(\begin{array}{c|c|c} \text{Image 1} & \dots & \text{Image } n \end{array} \right) \approx \left(\begin{array}{c|c|c|c|c} \text{Eigenface 1} & \text{Eigenface 2} & \text{Eigenface 3} & \text{Eigenface 4} & \text{Eigenface 5} \end{array} \right) \left(\begin{array}{c|c|c} \mathbf{z}_1 & \dots & \mathbf{z}_n \end{array} \right)$

Eigen-faces [Turk & Pentland 1991]

- d = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a face image
- \mathbf{x}_{ji} = intensity of the j -th pixel in image i

$$\mathbf{X}_{d \times n} \approx \mathbf{U}_{d \times k} \mathbf{Z}_{k \times n}$$

$\left(\begin{array}{c|c} \text{[Face 1]} & \dots & \text{[Face } n\text{]} \end{array} \right) \approx \left(\begin{array}{c|c|c|c|c} \text{[U1]} & \text{[U2]} & \text{[U3]} & \text{[U4]} & \text{[U5]} \end{array} \right) \left(\begin{array}{c|c|c} \text{[Z1]} & \dots & \text{[Zn]} \end{array} \right)$

Idea: \mathbf{z}_i more “meaningful” representation of i -th face than \mathbf{x}_i

Can use \mathbf{z}_i for nearest-neighbor classification

Eigen-faces [Turk & Pentland 1991]

- d = number of pixels
- Each $\mathbf{x}_i \in \mathbb{R}^d$ is a face image
- \mathbf{x}_{ji} = intensity of the j -th pixel in image i

$$\mathbf{X}_{d \times n} \approx \mathbf{U}_{d \times k} \mathbf{Z}_{k \times n}$$


The diagram illustrates the Eigenface decomposition. The matrix $\mathbf{X}_{d \times n}$ is shown as a row of face images. The matrix $\mathbf{U}_{d \times k}$ is shown as a row of five grayscale eigenface images. The matrix $\mathbf{Z}_{k \times n}$ is shown as a row of vectors $\mathbf{z}_1, \dots, \mathbf{z}_n$, each represented by a vertical line with a green dot at the top and a blue dot at the bottom.

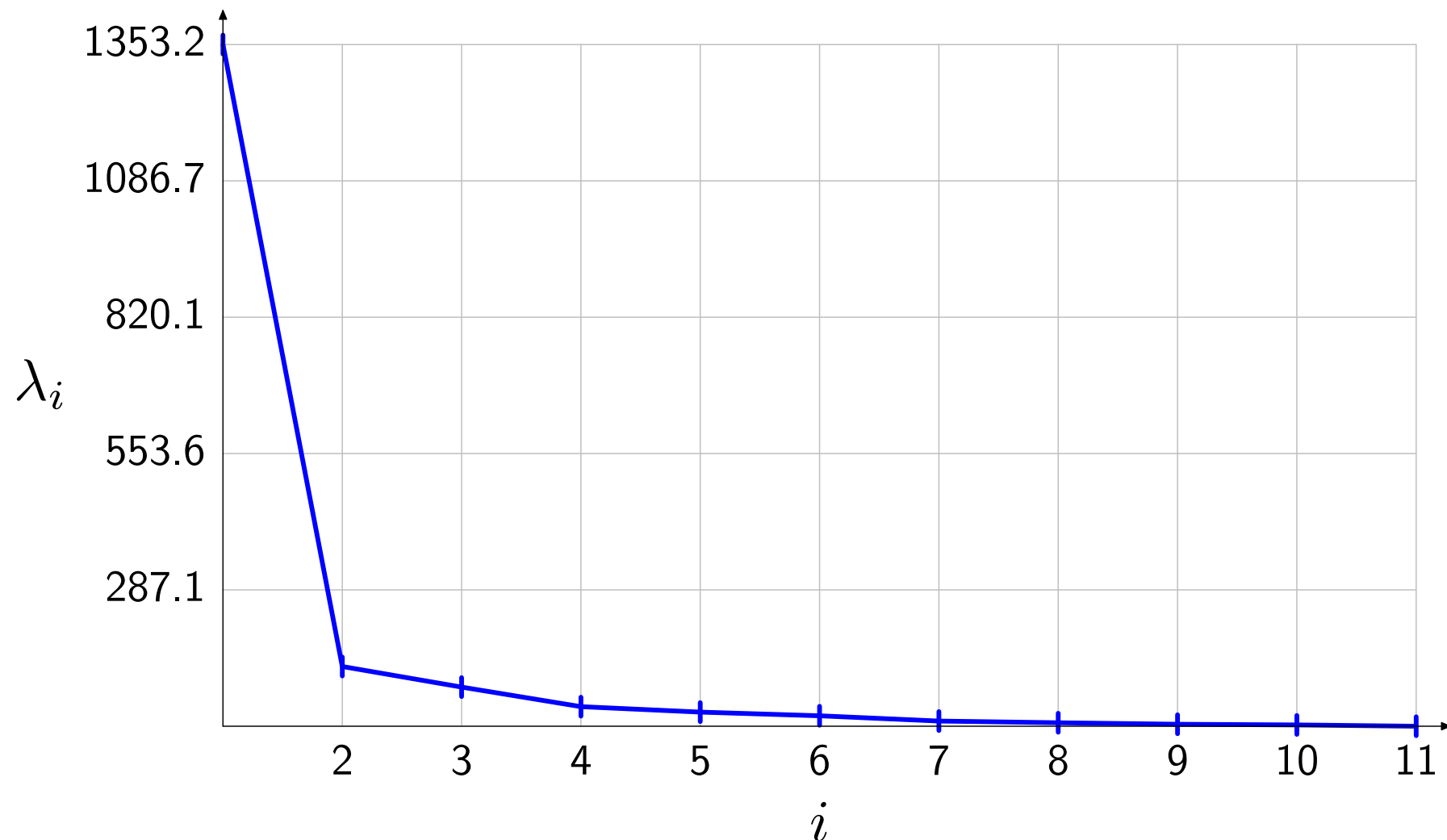
Idea: \mathbf{z}_i more “meaningful” representation of i -th face than \mathbf{x}_i

Can use \mathbf{z}_i for nearest-neighbor classification

Much faster: $O(dk + nk)$ time instead of $O(dn)$ when $n, d \gg k$

Aside: How many components?

- Magnitude of eigenvalues indicate fraction of variance captured.
- Eigenvalues on a face image dataset:



- Eigenvalues typically drop off sharply, so don't need that many.
- Of course variance isn't everything...

Wrapping up PCA

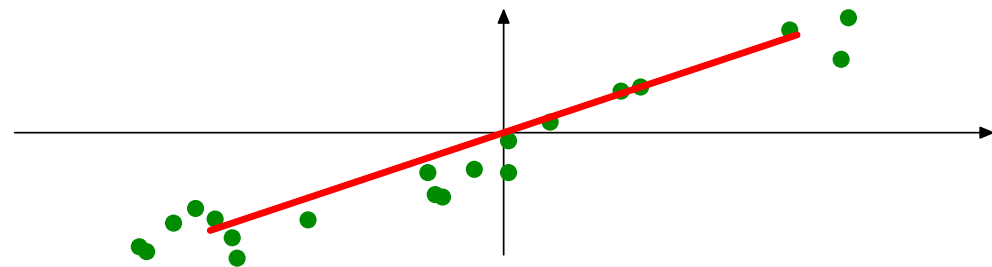
- PCA is a linear model for dimensionality reduction which finds a mapping to a lower dimensional space that maximizes variance

Wrapping up PCA

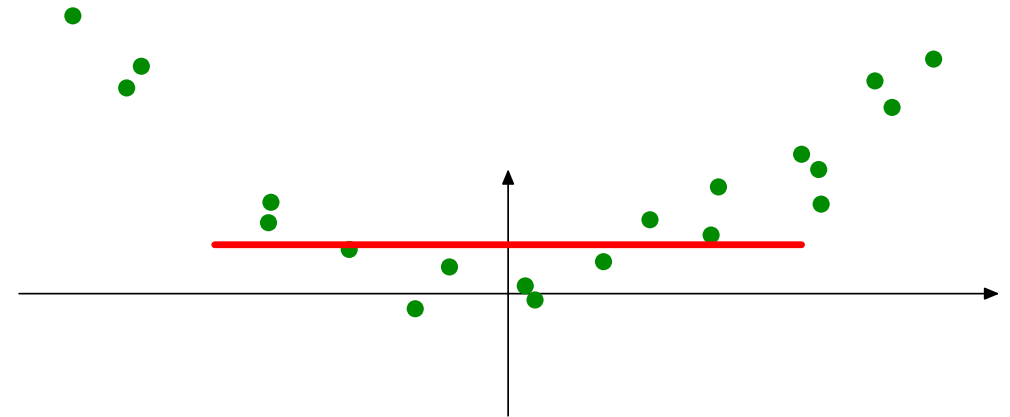
- PCA is a linear model for dimensionality reduction which finds a mapping to a lower dimensional space that maximizes variance
- We saw that this is equivalent to performing an eigendecomposition on the covariance matrix of \mathbf{X}

Non-linear dimensionality reduction

Limitations of Linearity

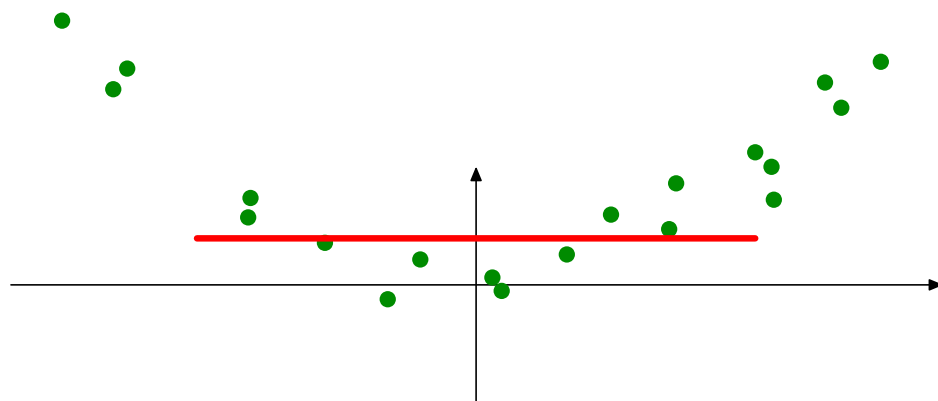


PCA is effective

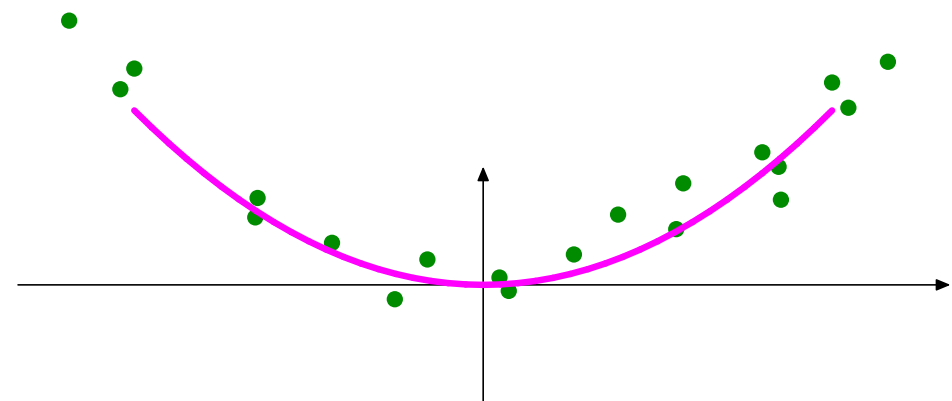


PCA is ineffective

Nonlinear PCA

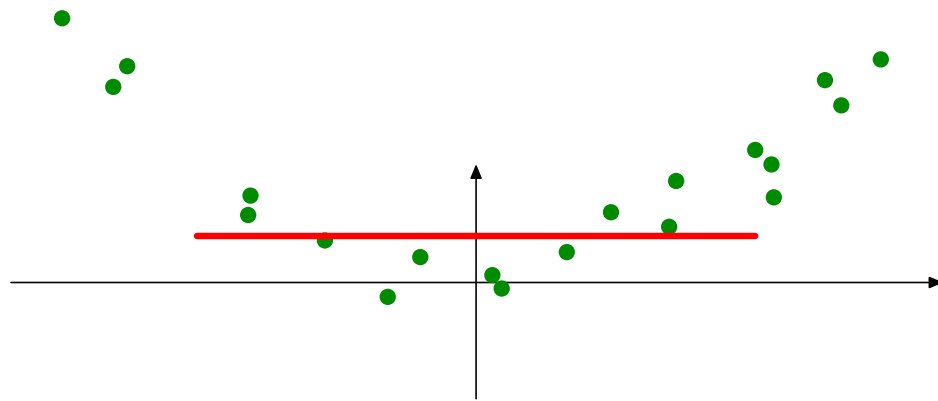


Broken solution

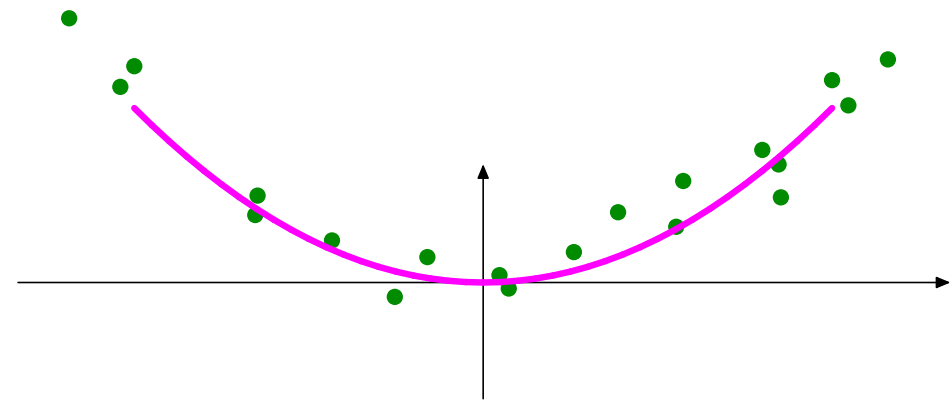


Desired solution

Nonlinear PCA



Broken solution



Desired solution

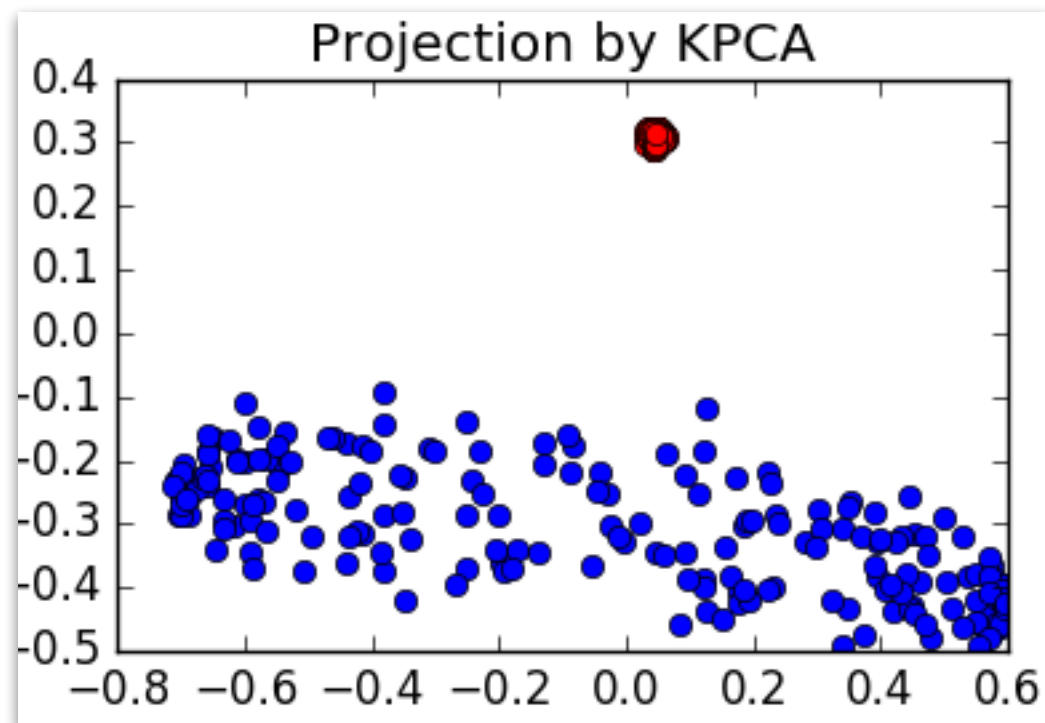
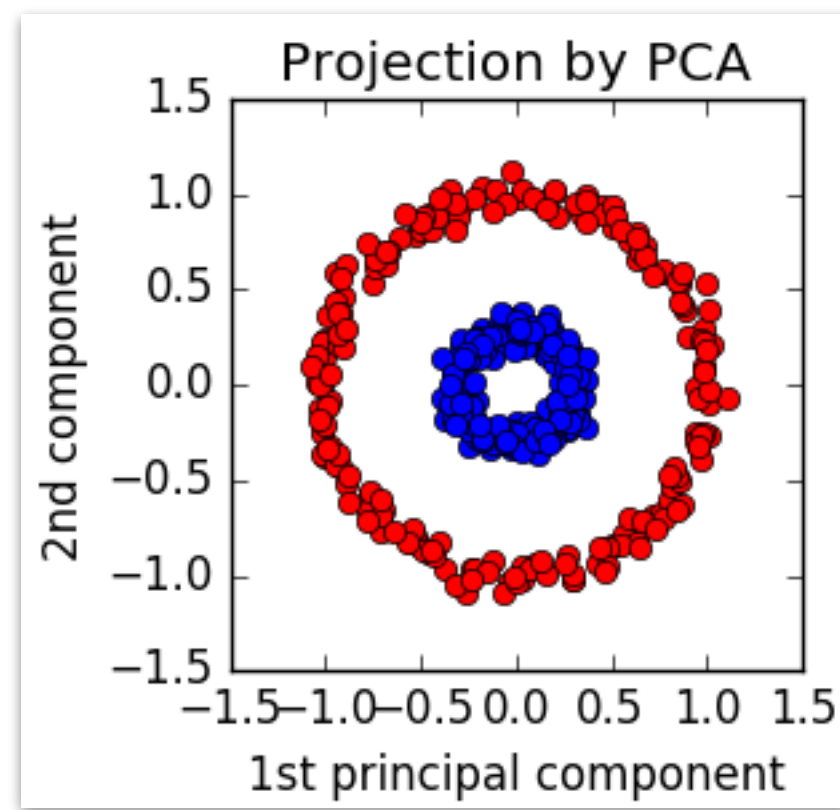
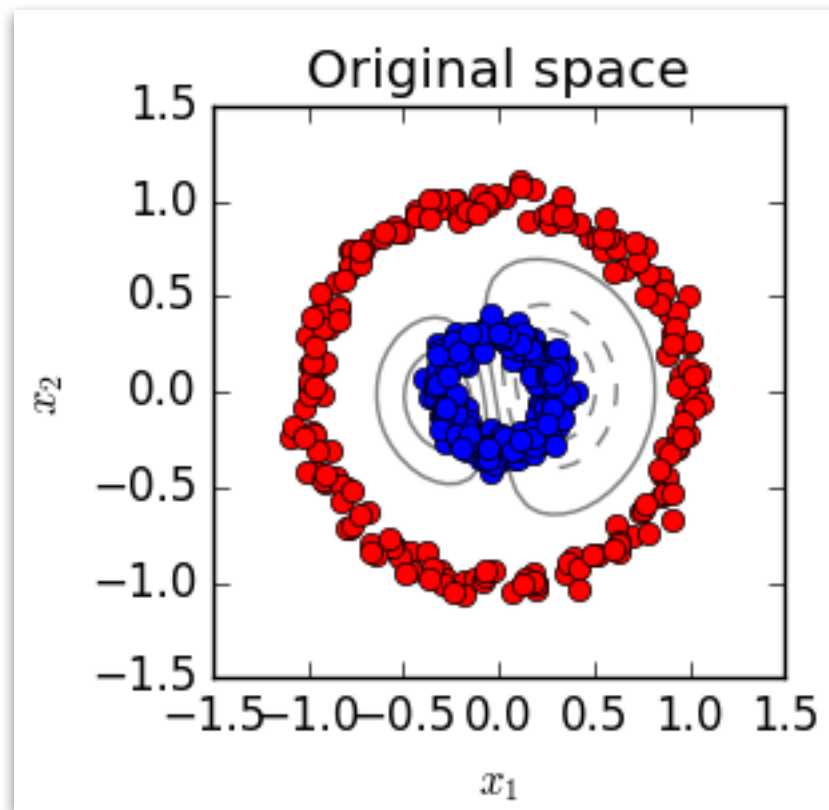
Idea: Use kernels

Linear dimensionality reduction in $\phi(\mathbf{x})$ space



Nonlinear dimensionality reduction in \mathbf{x} space

Kernel PCA



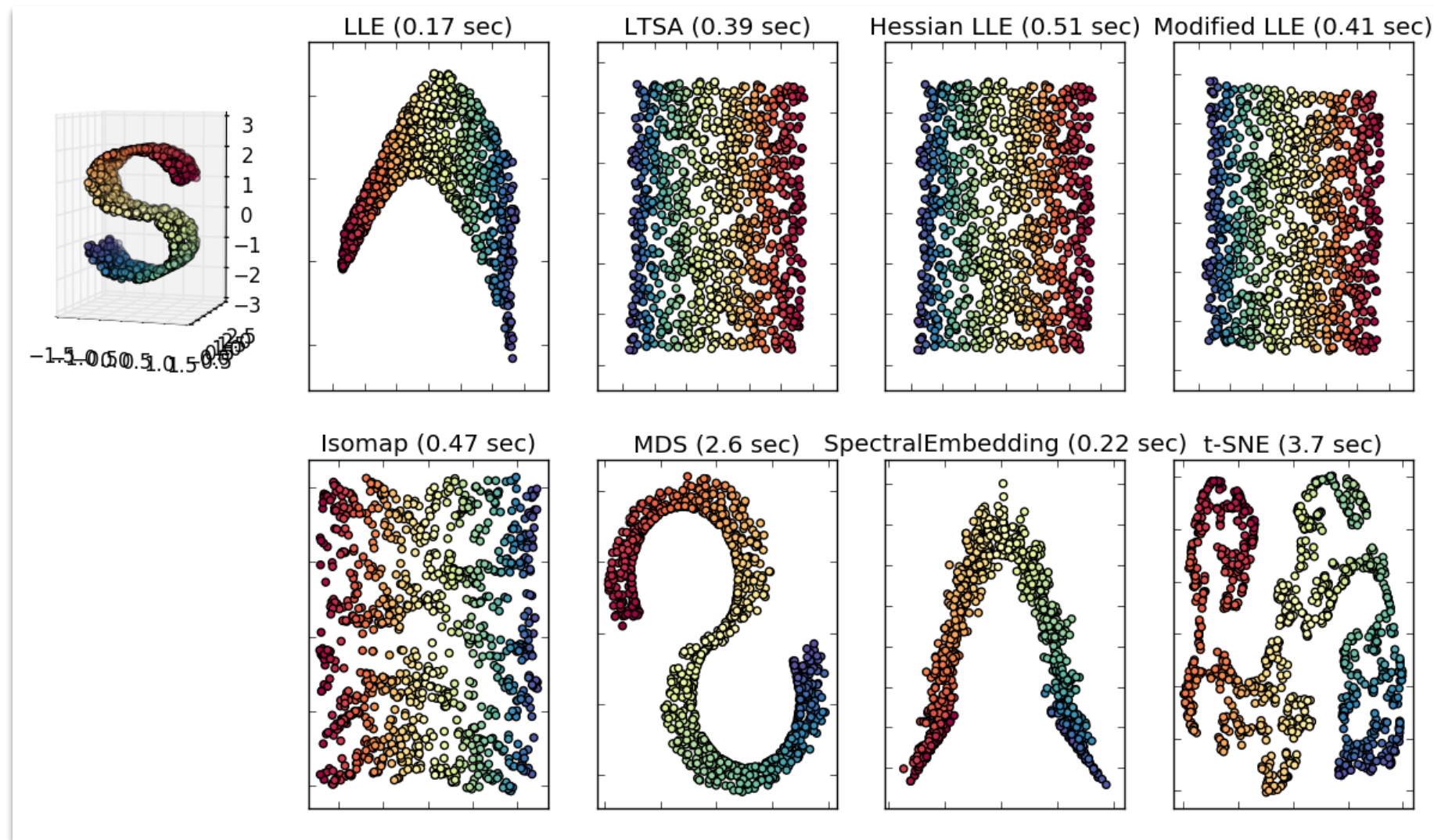
Alternatively: *t-SNE*!

Stochastic Neighbor Embeddings



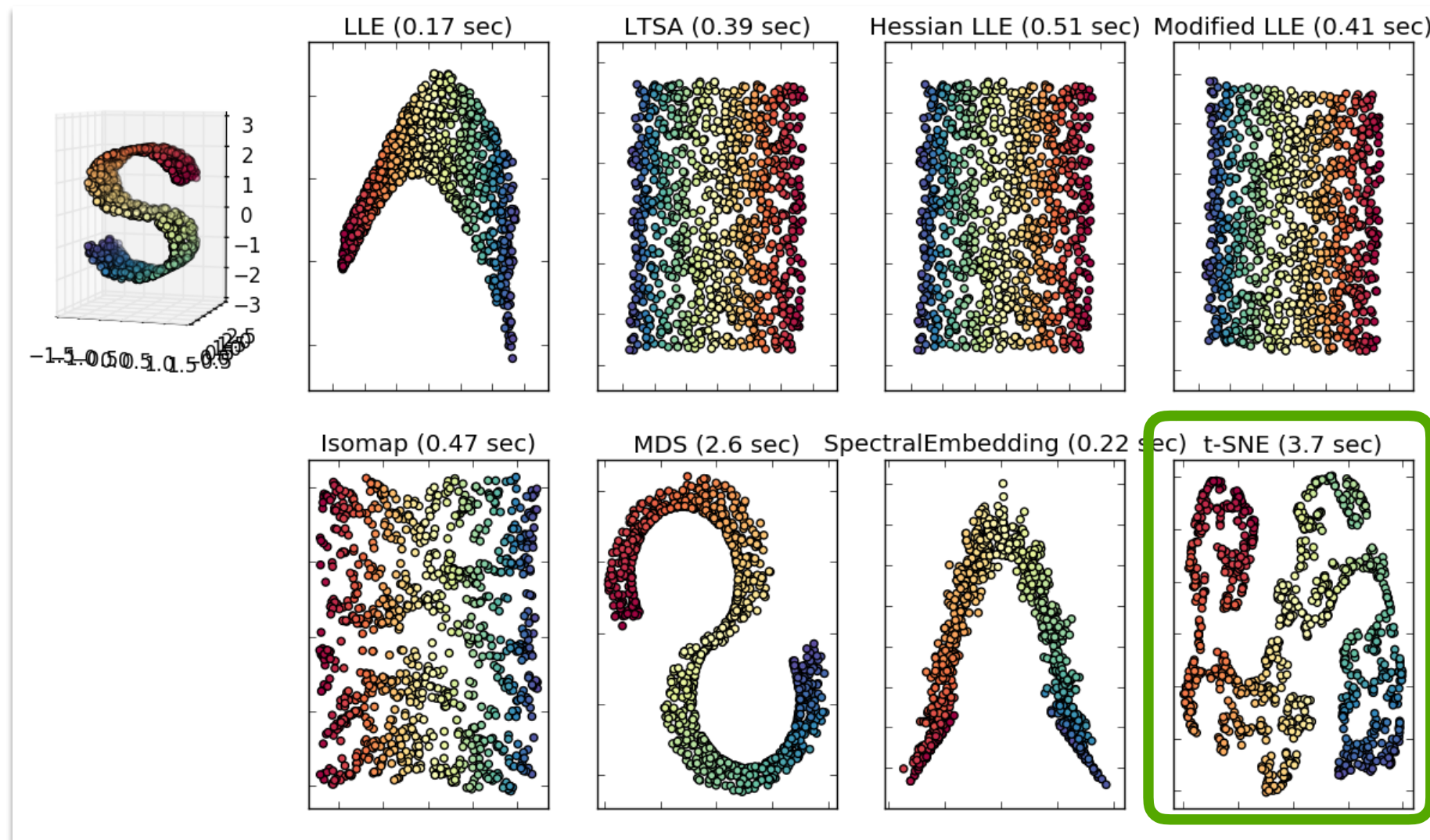
Borrowing from:
Laurens van der Maaten
(Delft -> Facebook AI)

Manifold learning



Idea: Perform a *non-linear* dimensionality reduction in a manner that preserves proximity (but not distances)

Manifold learning



Visualizing data using t-SNE

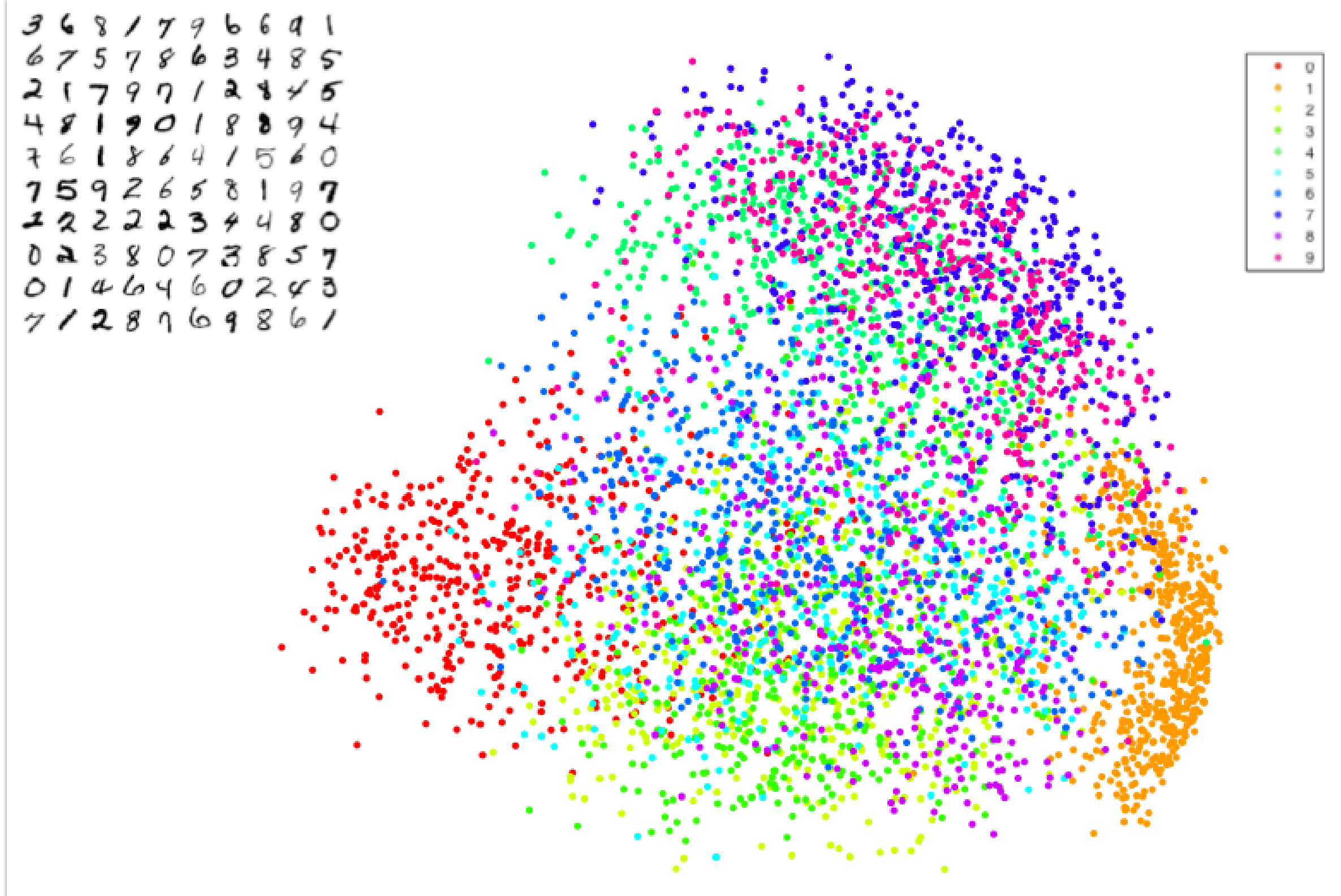
[L Maaten, G Hinton](#) - Journal of machine learning research, 2008 - jmlr.org

[+ Paperpile](#)

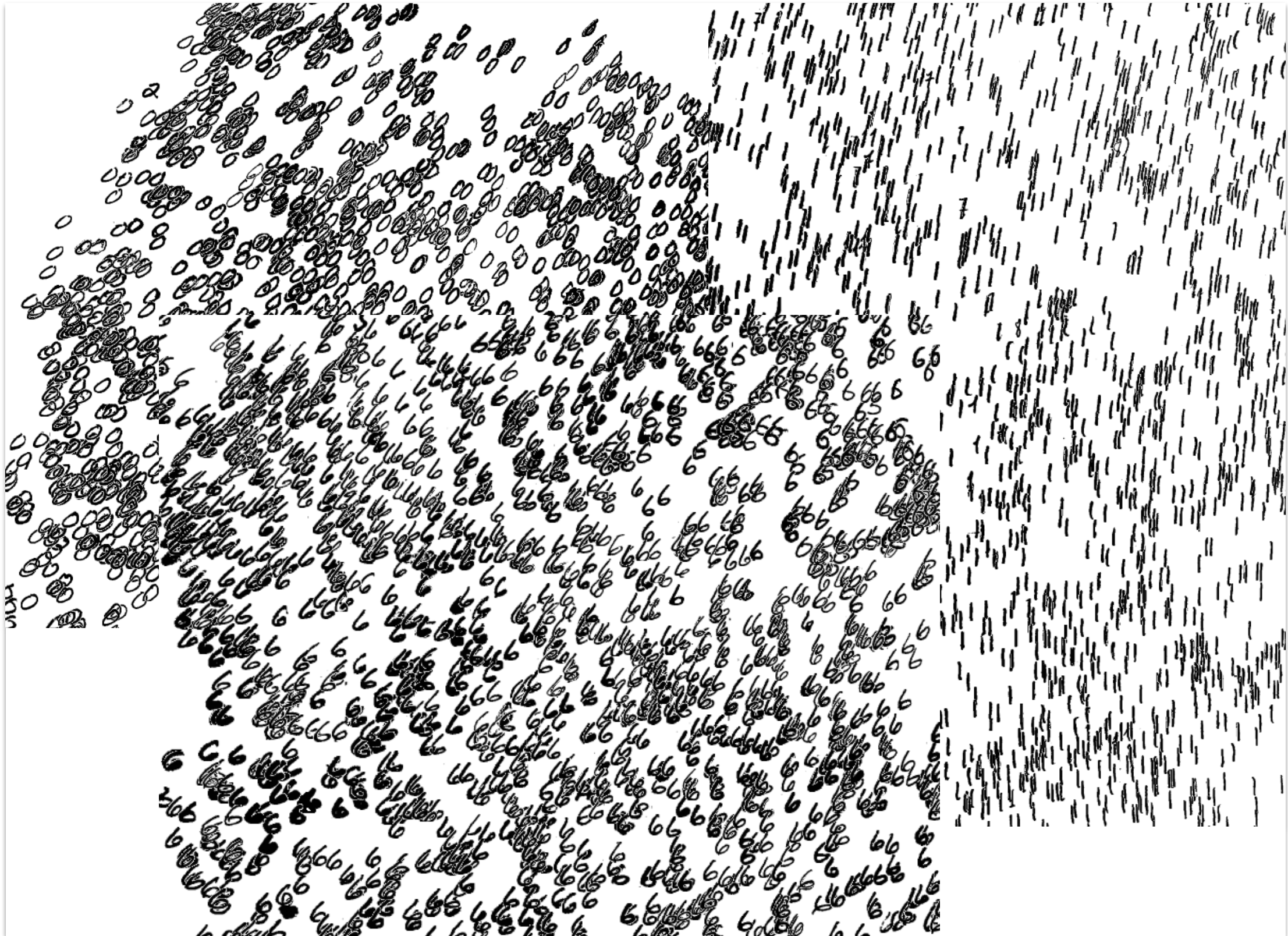
We present a new technique called "t-SNE" that visualizes high-dimensional data by giving each datapoint a location in a two or three-dimensional map. The technique is a variation of Stochastic Neighbor Embedding (Hinton and Roweis, 2002) that is much easier to optimize ...

☆ [Cited by 11621](#) [Related articles](#) [All 39 versions](#) [Import into BibTeX](#) [»»](#)

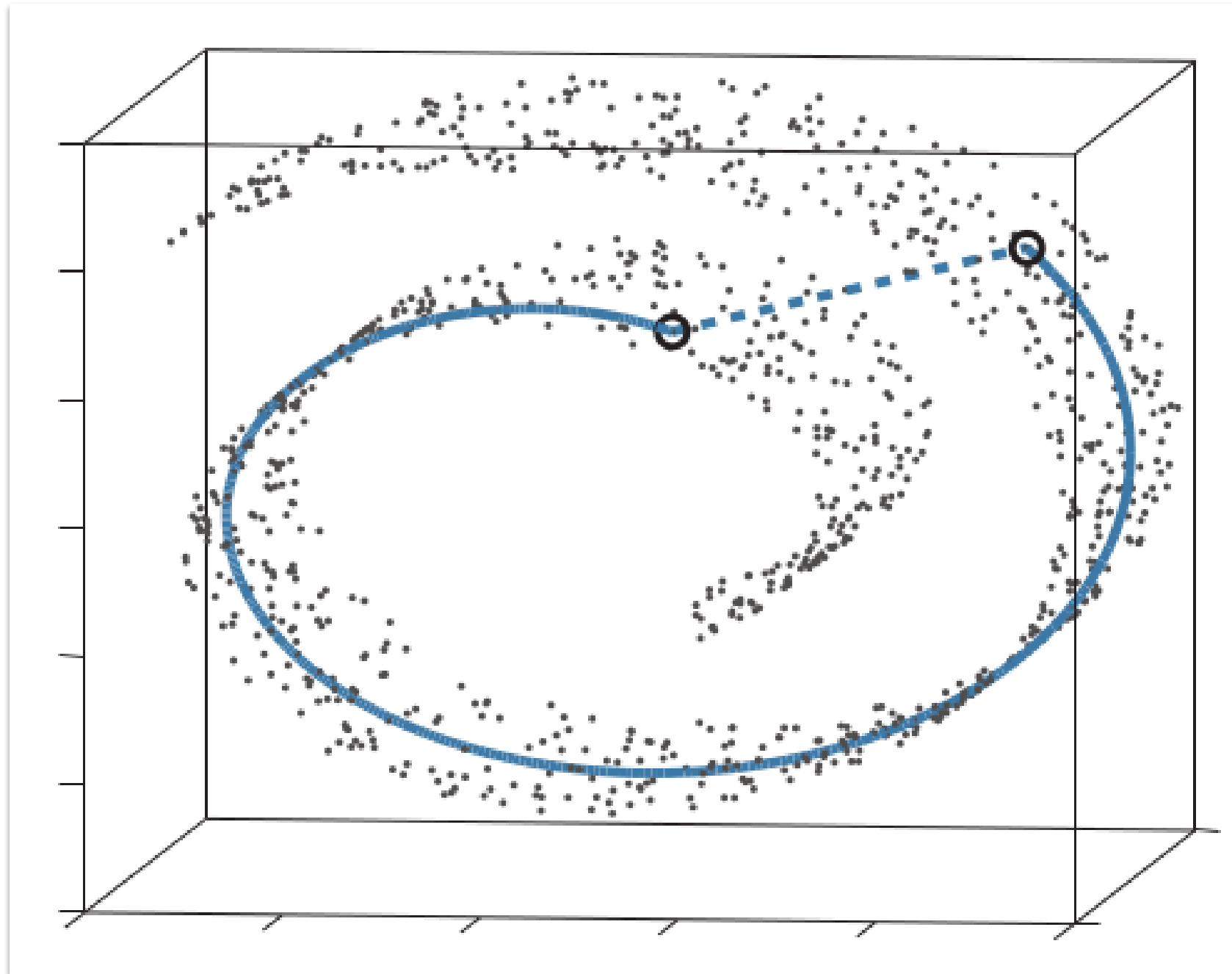
PCA on MNIST digits



t-SNE on MNIST Digits

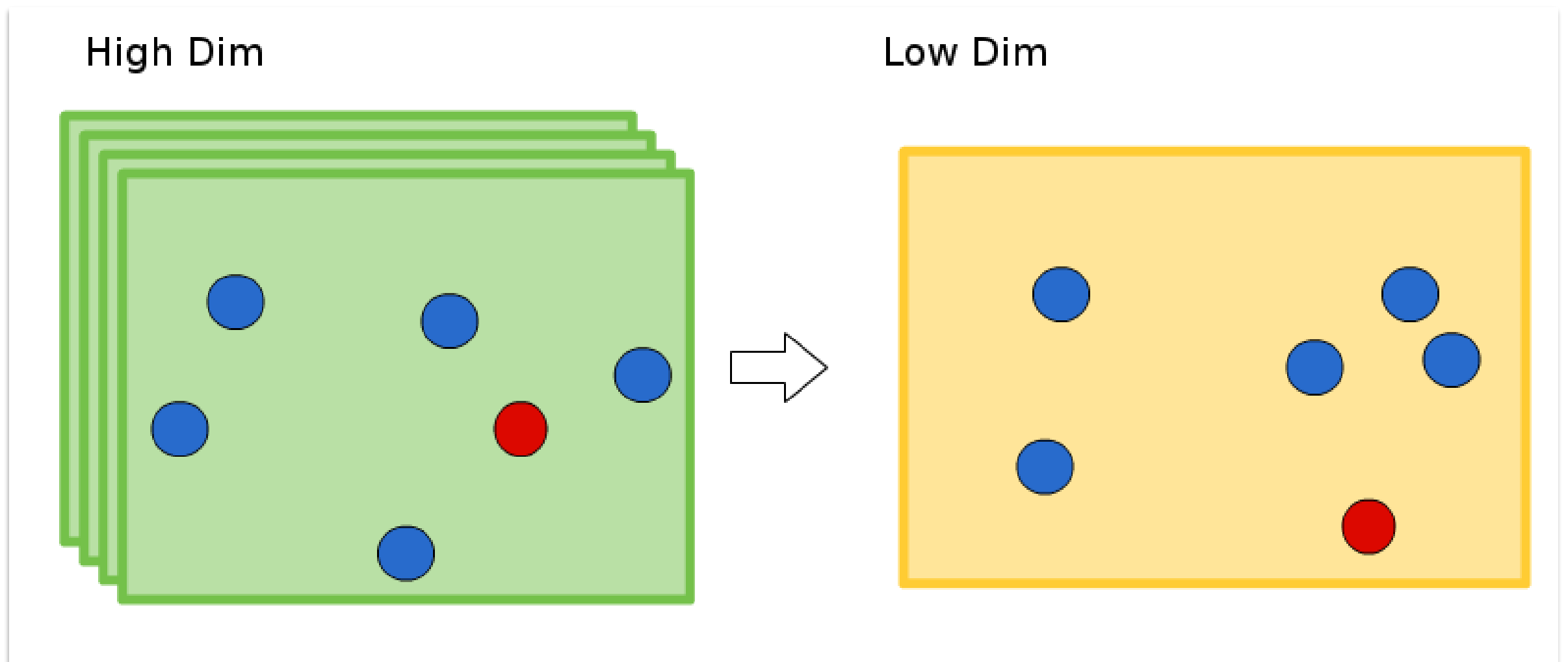


Swiss roll



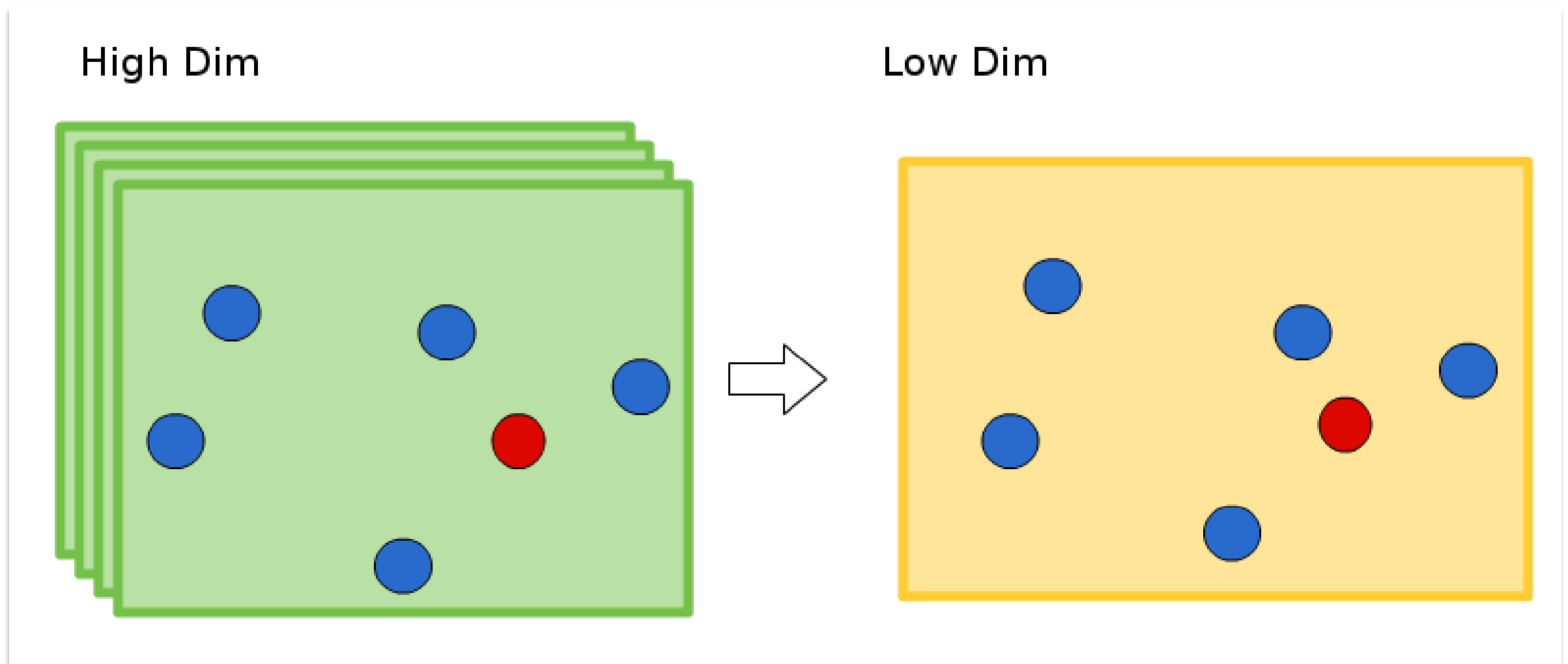
Euclidean distance is not always
a good notion of proximity

Non-linear projection



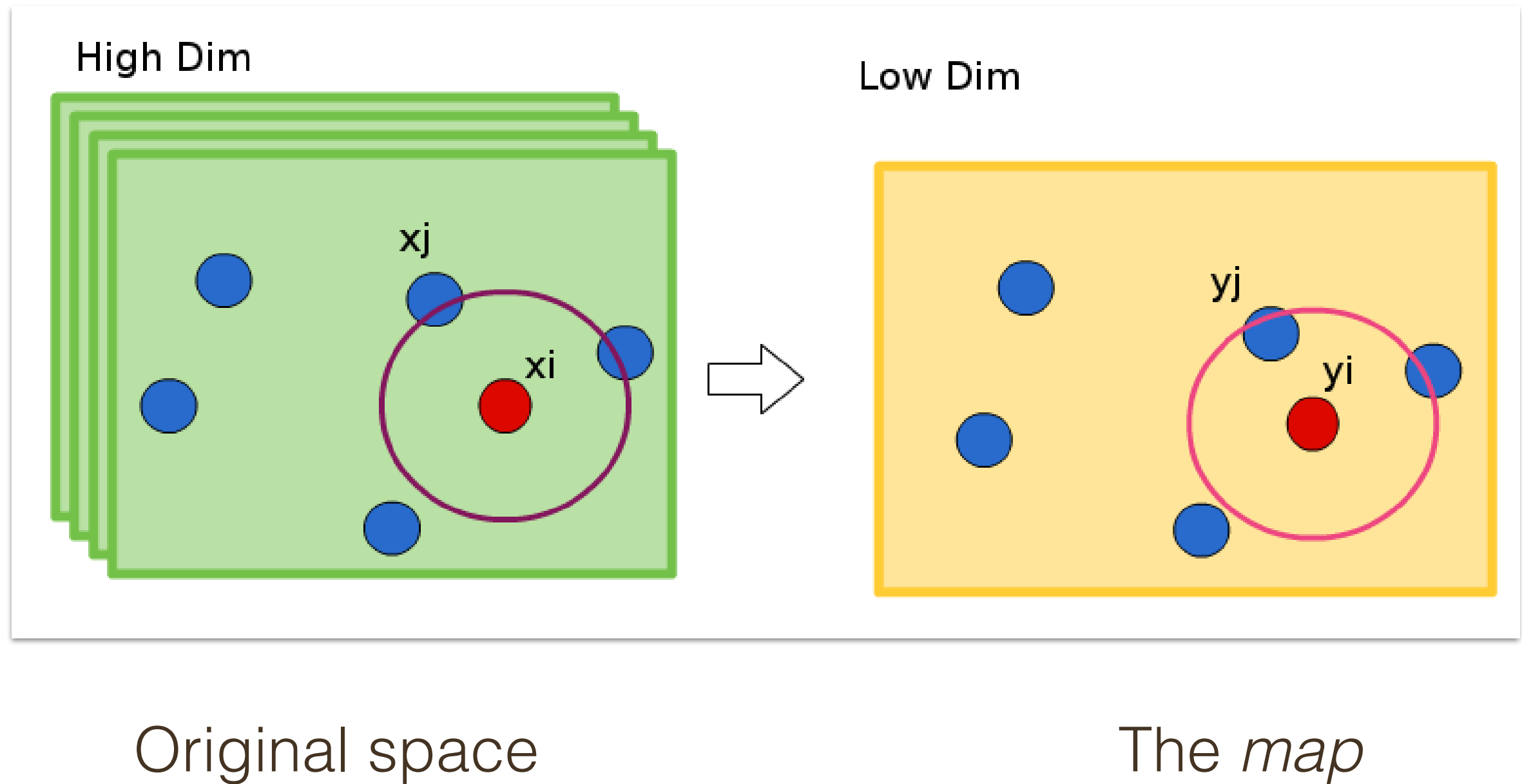
Bad projection: relative position to neighbors changes

Non-linear projection



Intuition: Want to preserve ***local*** neighborhood

Stochastic Neighbor Embedding



SNE to t-SNE
(on board)

t-SNE: SNE with a t-Distribution

Similarity in High Dimension

$$p_{ij} = \frac{\exp(-||x_i - x_j||^2 / 2\sigma^2)}{\sum_{k \neq i} \exp(-||x_i - x_k||^2 / 2\sigma^2)}$$

Similarity in Low Dimension

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

Gradient

$$\frac{\partial C}{\partial y_i} = 4 \sum_{j \neq i} (p_{ij} - q_{ij})(1 + ||y_i - y_j||^2)^{-1}(y_i - y_j)$$

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

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

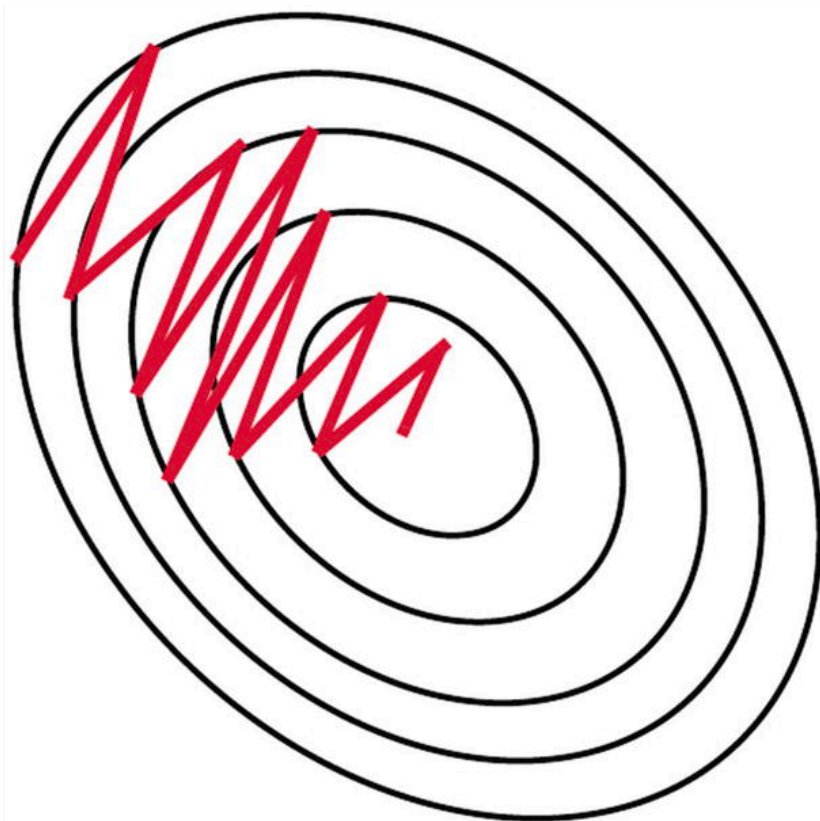
$$\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta \mathcal{C}}{\delta \mathcal{Y}} + \alpha(t) (\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)})$$

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

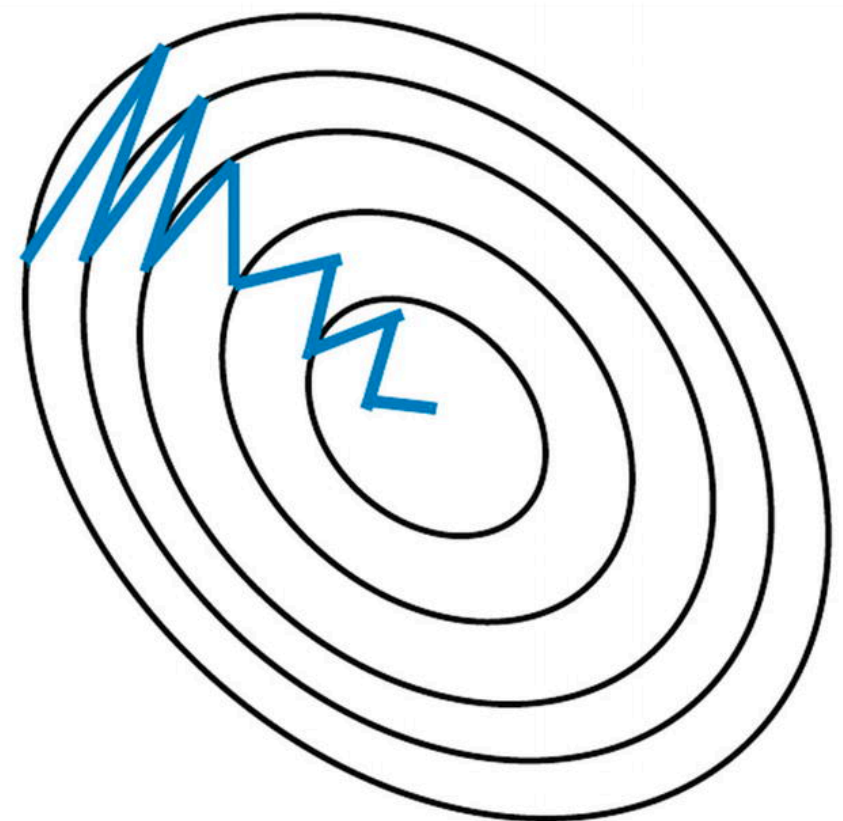
Regular gradient descent

“momentum”

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



Stochastic Gradient
Descent **without**
Momentum

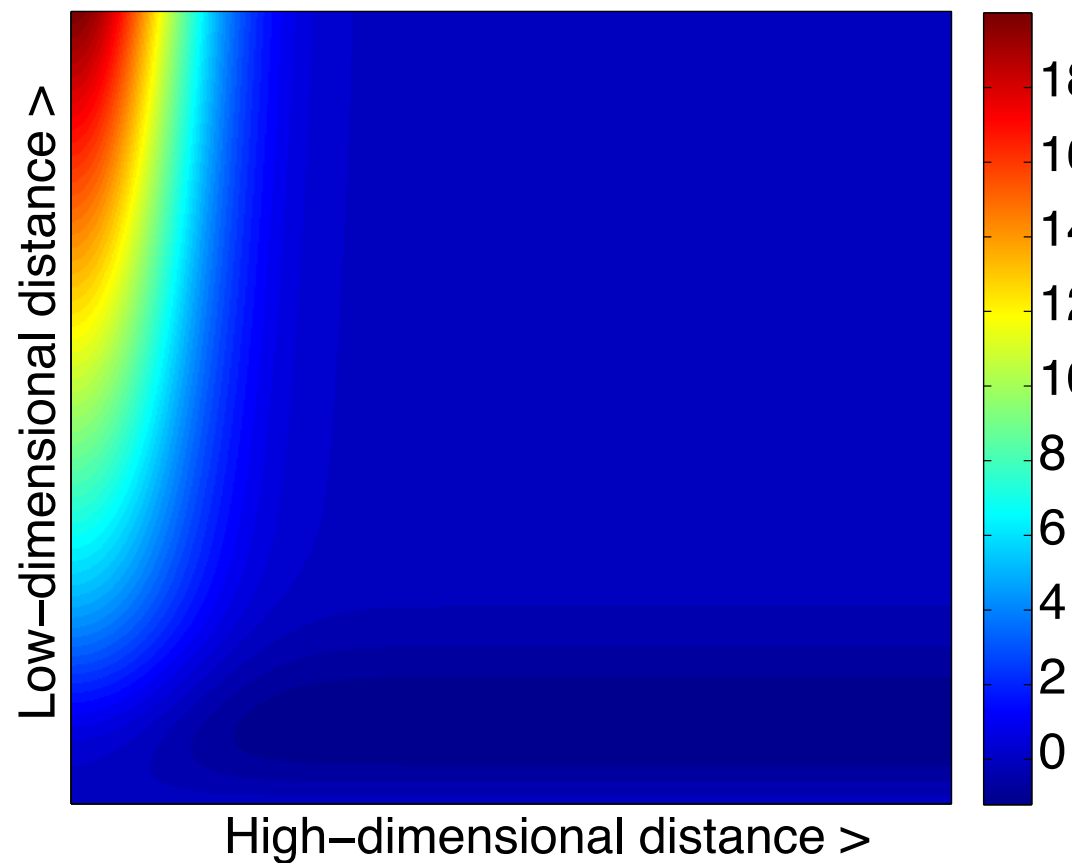


Stochastic Gradient
Descent **with**
Momentum

Figure credit: Bisong, 2019

What's magical about t ?

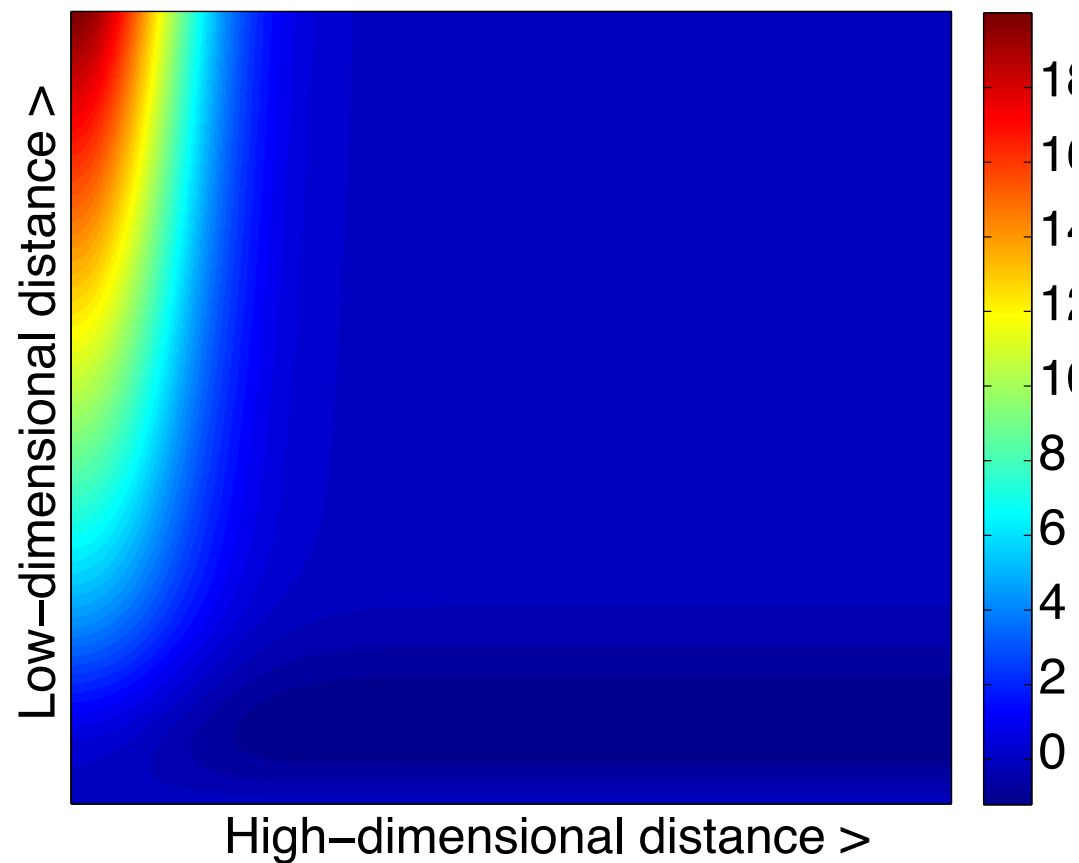
Basically, the gradient has nice properties



(a) Gradient of SNE.

What's magical about t ?

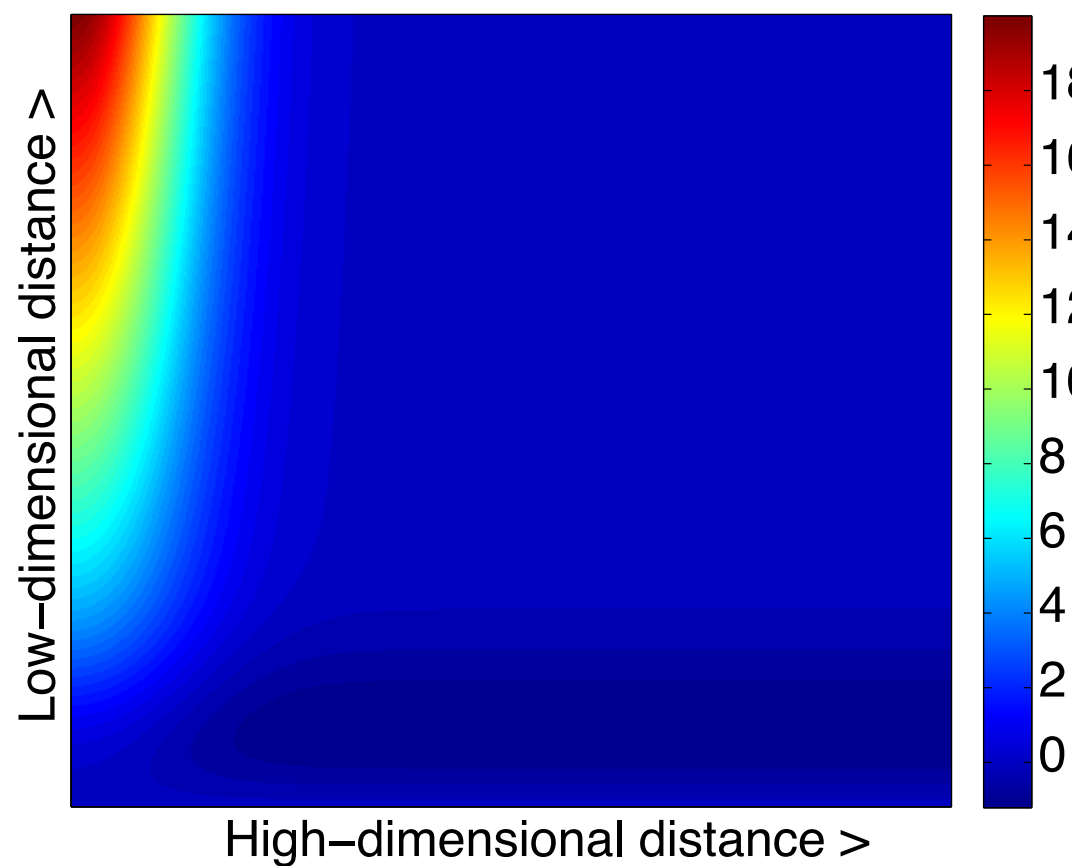
Basically, the gradient has nice properties



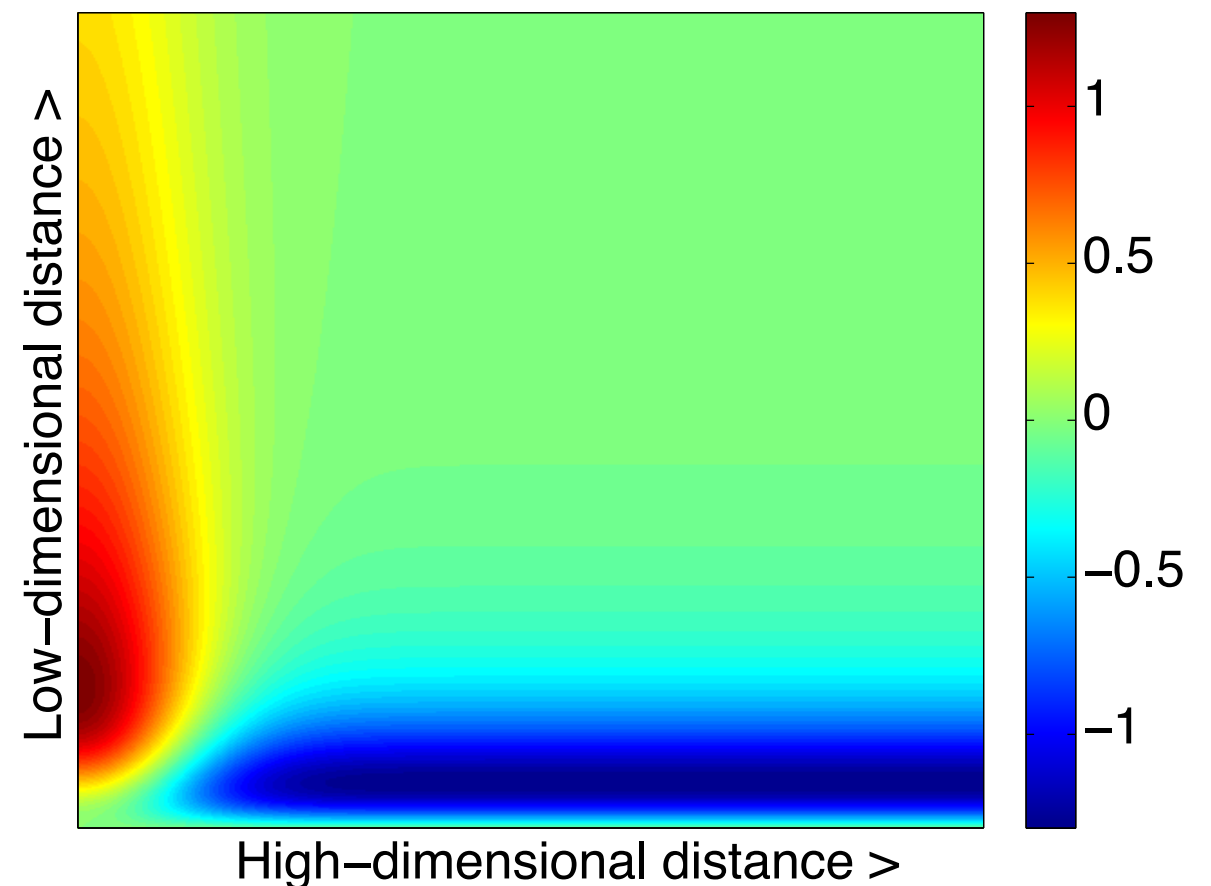
(a) Gradient of SNE.

Positive gradient \rightarrow “attraction” between points

What's magical about t ?



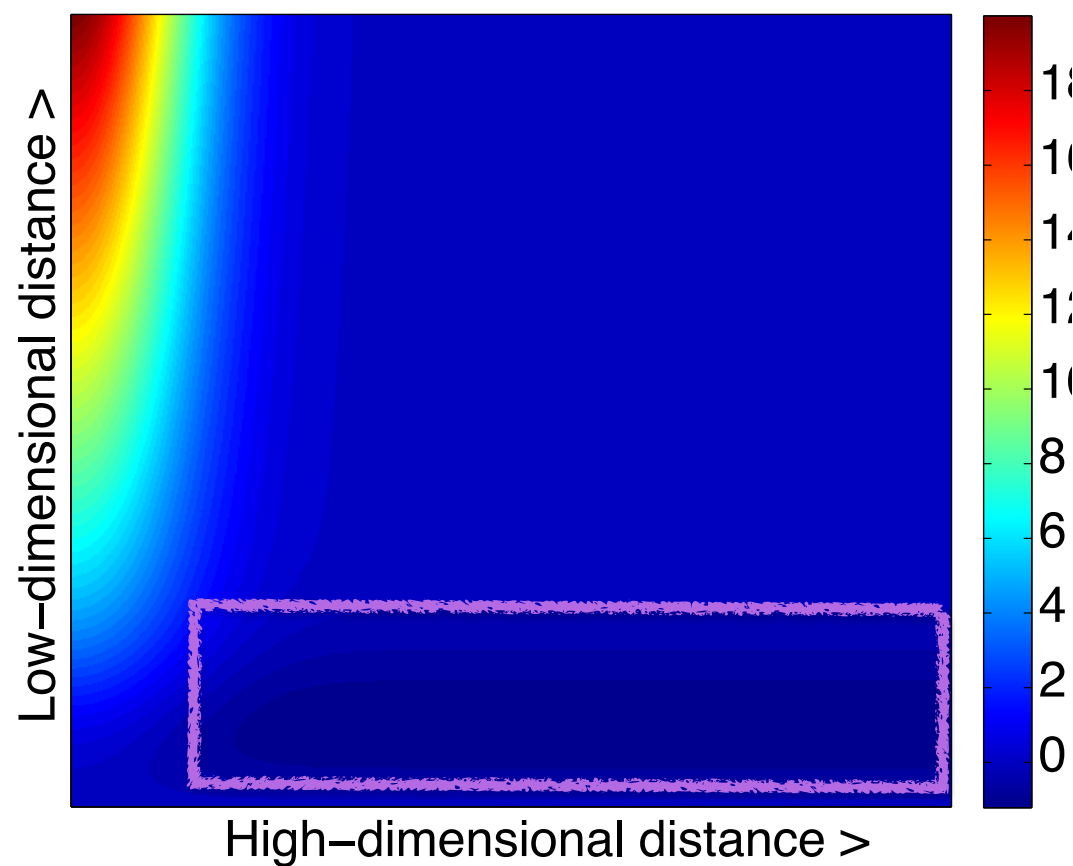
(a) Gradient of SNE.



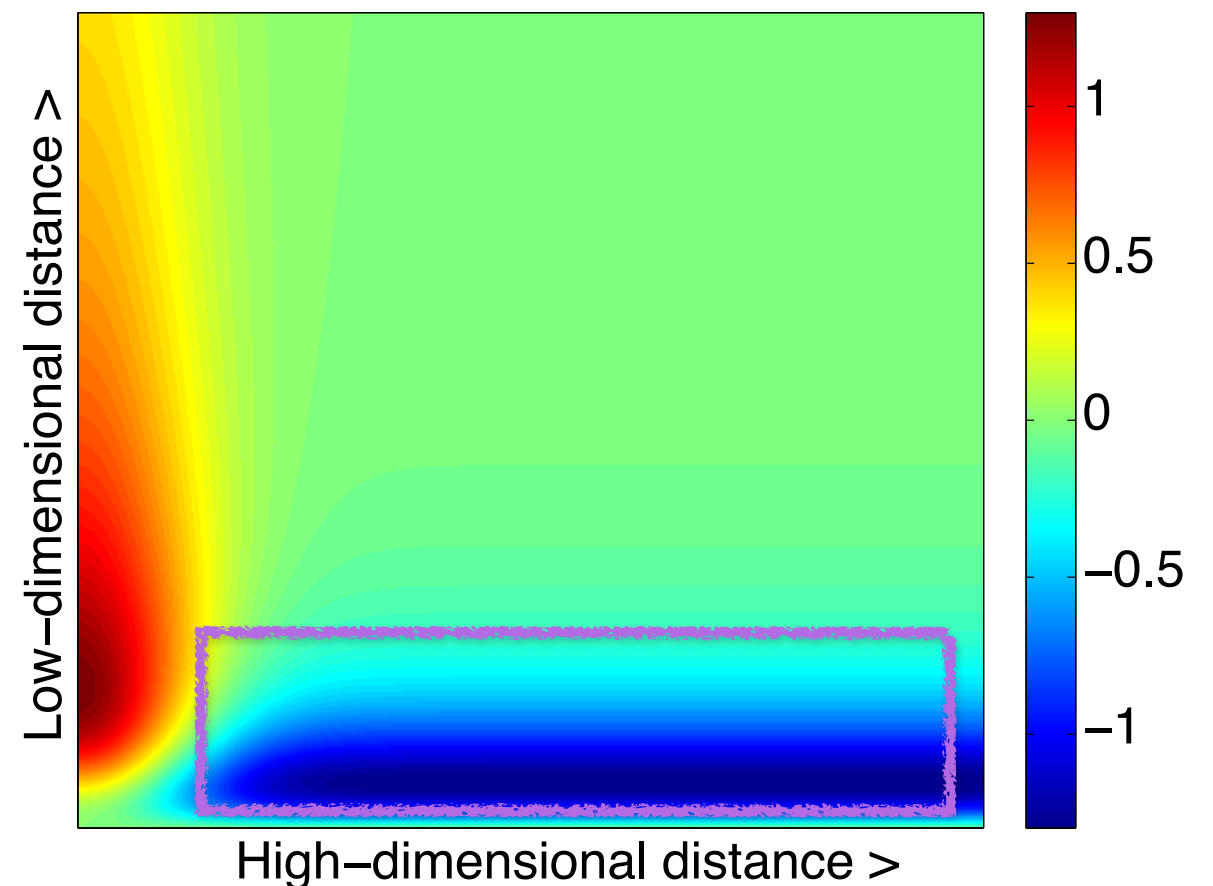
(c) Gradient of t-SNE.

Positive gradient \rightarrow “attraction” between points

What's magical about t ?



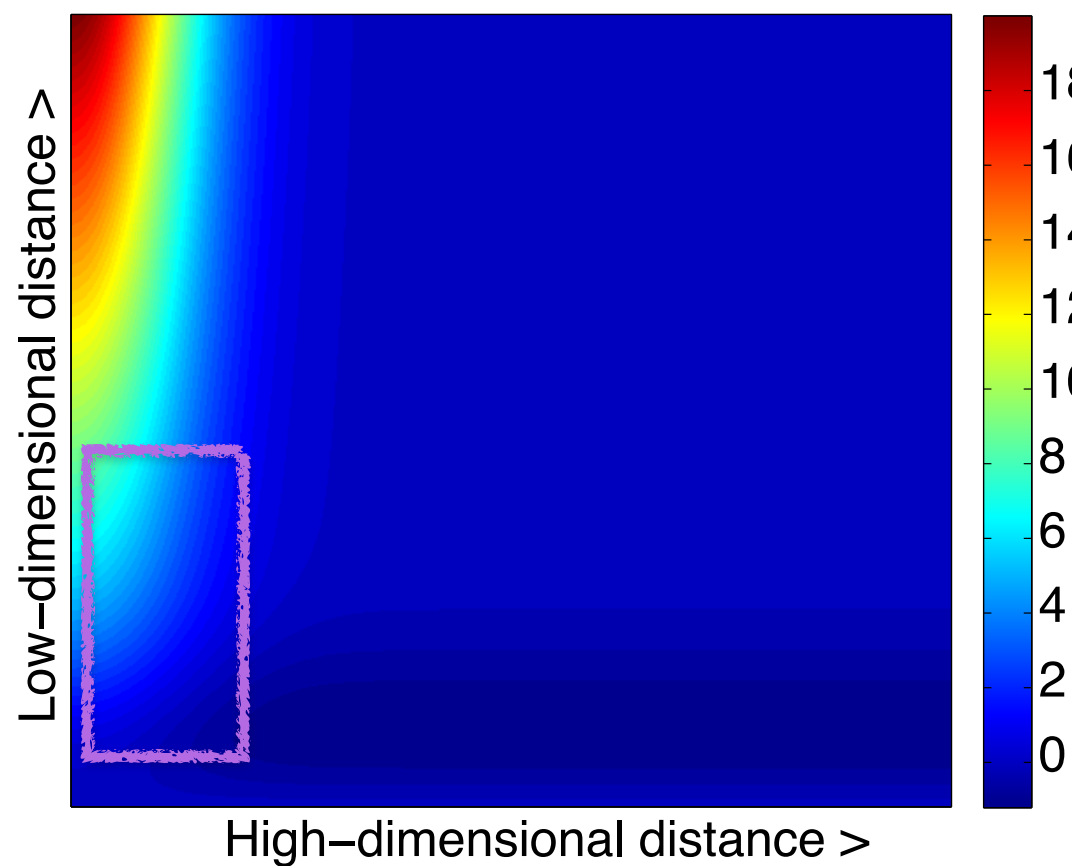
(a) Gradient of SNE.



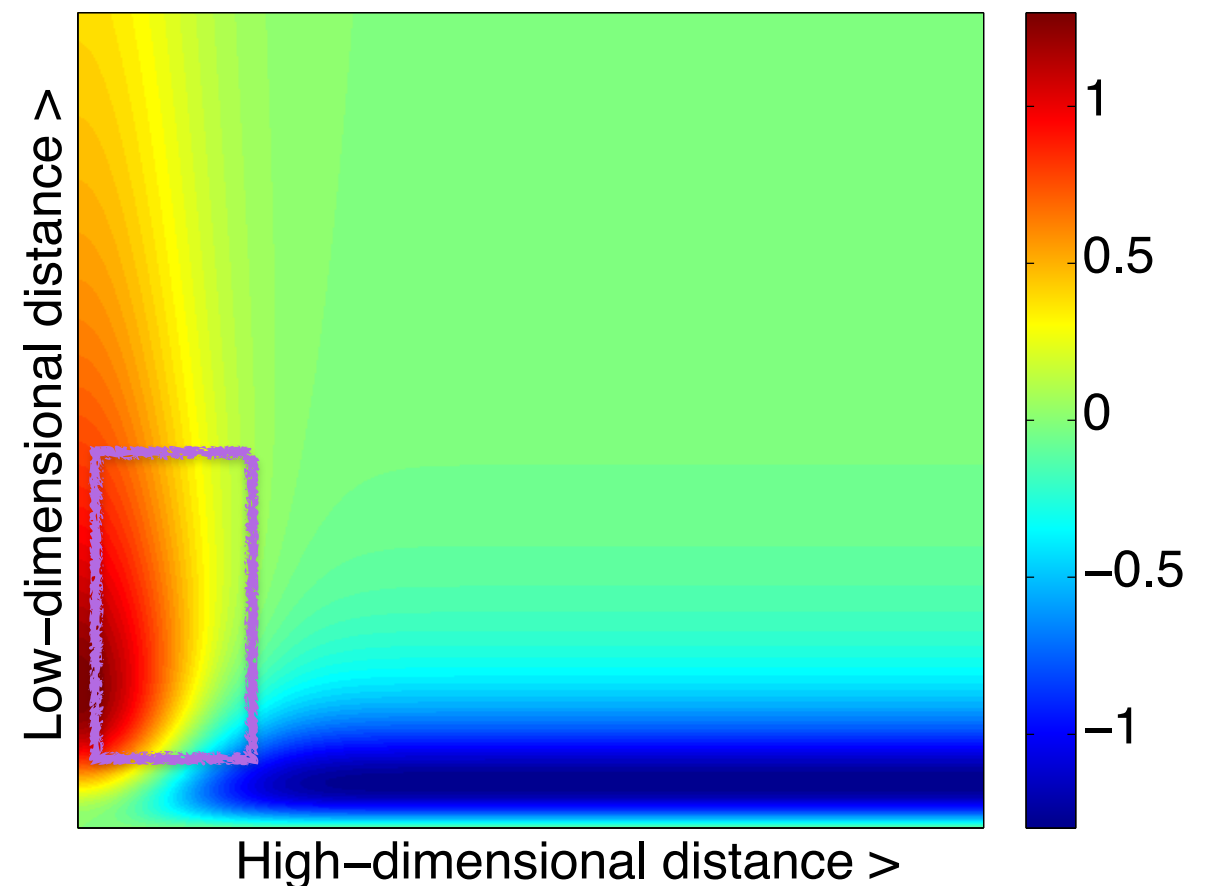
(c) Gradient of t-SNE.

t-SNE *repels* points in low dim space that are different in the high dim space

What's magical about t ?



(a) Gradient of SNE.



(c) Gradient of t-SNE.

Also strongly attracts points nearby in high dim space

Let's see some code

Another perspective:
Auto-encoders

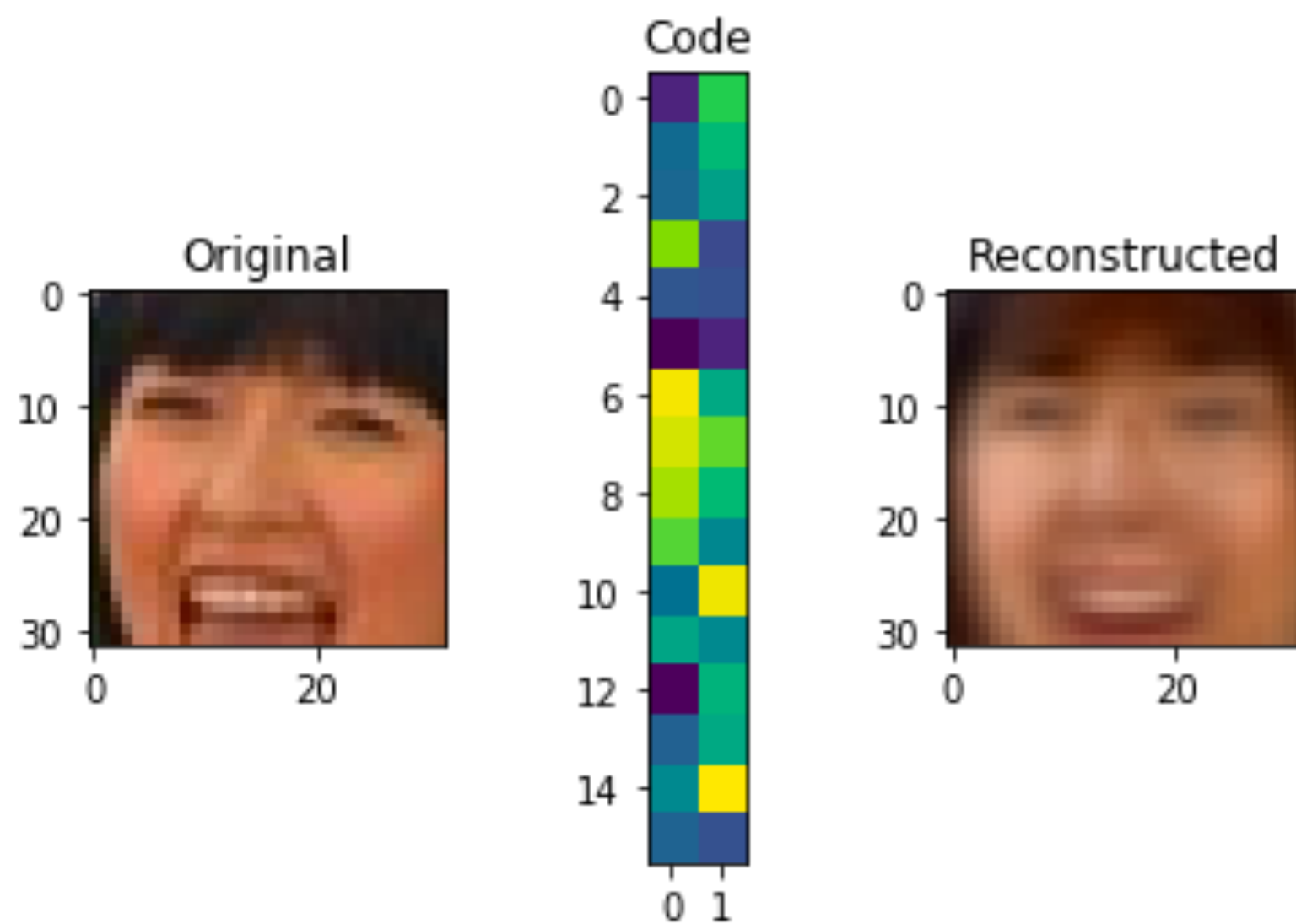


Figure credit: <https://stackabuse.com/autoencoders-for-image-reconstruction-in-python-and-keras/>