

- kernel trick: replace G with a kernel matrix $K_{i,i'} = \underbrace{K(X_i, X_{i'})}_{\text{kernel function}}$

and compute eigen decomposition of K

- catch: PCA requires the data to be centered, but the K -matrix is not centered: $K_{i,i'} = \varphi(X_i) \cdot \varphi(X_{i'})^T$ for some $\varphi(x)$
 $\varphi(x)$ is implicitly defined by $K(X_i, X_{i'})$ and generally not centered
- Q: how can we center K without computing $\bar{\varphi}(x) = \frac{1}{N} \sum_i \varphi(X_i)$

A: we can compute a centered Gram matrix (kernel matrix) from uncentered data!

let X' be uncentered data, $X \rightarrow X' - \bar{X}$ $X = X' - \frac{1}{N} \cdot \bar{X}$
 G' uncentered Gram matrix $G' = X' X'^T$ \leftarrow column vector of all 1's

G centered Gram matrix $G = X X^T = (X' - \frac{1}{N} \mathbb{1} \bar{X}') (X' - \frac{1}{N} \mathbb{1} \bar{X}')^T$

$$G = \underbrace{X' X'^T}_{= G'} - \frac{1}{N} \mathbb{1} \bar{X}' X'^T - X' \left(\frac{1}{N} \bar{X}' \right)^T + \frac{1}{N} \mathbb{1} \bar{X}' \bar{X}'^T \frac{1}{N} \mathbb{1}^T$$

$$= G' - \frac{1}{N} \mathbb{1} \mathbb{1}^T \underbrace{X' X'^T}_{= G'} - X' \underbrace{\left(\frac{1}{N} \bar{X}' \right)^T}_{\frac{1}{N} X'^T \mathbb{1}} + \frac{1}{N} \mathbb{1} \mathbb{1}^T X' \left(\frac{1}{N} X'^T \mathbb{1} \right) \frac{1}{N} \mathbb{1}^T$$

$$= G' - \underbrace{\mathbb{1} \left(\frac{1}{N} \mathbb{1}^T G' \right)}_{\bar{G}' \text{ row vector of column averages}} - \underbrace{\left(\frac{1}{N} G' \mathbb{1} \right) \mathbb{1}^T}_{\bar{G}'^T} + \mathbb{1} \underbrace{\left(\frac{1}{N^2} \mathbb{1}^T G' \mathbb{1} \right) \mathbb{1}^T}_{\gamma: \text{average value of } G'}$$

Post-centering of Gram matrix:

- (1) compute $G' = X' X'^T$ uncentered Gram matrix
- (2) compute row matrix of column averages $\bar{G}' = \frac{1}{n} \mathbb{1}^T G'$ and total average
and total average $\gamma = \frac{1}{n^2} \mathbb{1}^T G' \mathbb{1}$

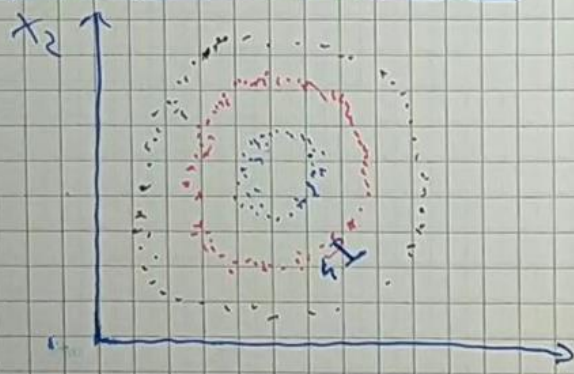
- (3) compute centered Gram matrix $G = G' - \mathbb{1} \bar{G}' - \bar{G}'^T \mathbb{1}^T + \gamma \mathbb{1} \mathbb{1}^T$

This also works when $G'_{ij} = k(x_i, x_j)$ with kernel trick

⇒ kernel PCA:

- (1) define G' from $k(x_i, x_j)$
- (2) center: $G = G' - \dots$ (as above)
- (3) define new features $z_j = u_j \lambda_j^{1/2}$ with $G u = \lambda u$ the eigen-decomp.

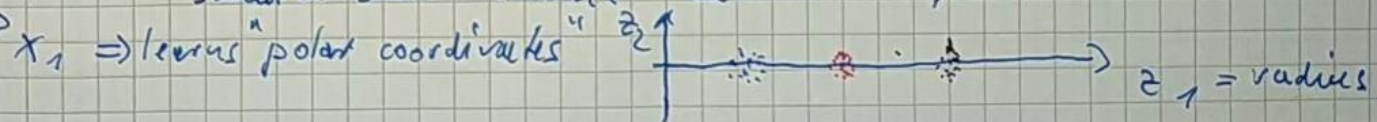
why is this case fun?



PCA has no effect: data form a circle \Rightarrow features

x_1 and x_2 are already uncorrelated, eigenvectors
 \checkmark of I are unit vectors, $z = x$

kernel PCA with squared exponential kernel with correct
bandwidth γ larger than dist within rings, and
smaller than distance between rings



Dimension reduction: compute new features $z = \varphi(X)$

$$\dim(z_i) \stackrel{\ll}{<} \dim(x_i), \quad \varphi \begin{cases} \text{linear (PCA)} \\ \text{non-linear (t-SNE)} \end{cases}$$

- purposes:
 - use z as input of some other ML algorithm,
 - instead of $X \Rightarrow$ may run faster, because lower dimension,
 - may be more accurate, because z contains the essence of X but not its noise

- use z for visualization: $\dim(z) = 2$ or 3 , draw as picture
 \Rightarrow the essential behavior of X may still be visible (e.g. clusters)

• further non-linear algorithms: Local Linear Embedding (LLE)

- idea: - find for each i the k nearest neighbors according to distance $d(x_i, x_{i'})$

- express x_i as a linear combination of its neighbors: chosen by designer

$$x_i = \sum_{i' \in \text{Nei}(x_i)} w_{i'} \cdot x_{i'}, \quad \text{with } \sum_{i'} w_{i'} = 1$$

- choose new coordinates z_i such that the relation with neighbors is approximately preserved:

$$z_i \approx \sum_{i' \in \text{Nei}(x_i)} w_{i'} z_{i'}$$

\Rightarrow effect: local relationships between instances are preserved, long-distance relations may change a lot.

- optimization problem: $\hat{w}_{i'}^{(i)} = \arg \min_w \|x_i - w X\|^2$ s.t. $w_{i'}^{(i)} = 0$ for $i' \notin \text{Nei}(x_i)$
 and $\sum_{i'} w_{i'}^{(i)} = 1$

simplify by defining $\tilde{X}^{(i)}$: matrix of only neighbors of X_i

$$\sum_{i'} w_{i'} = 1 \Leftrightarrow w \cdot \mathbb{1} = 1$$

$$G^{(i)} = \left(\tilde{X}^{(i)} - \mathbb{1} X_i \right) \left(\tilde{X}^{(i)} - \mathbb{1} X_i \right)^T$$

centering of $\tilde{X}^{(i)}$ around X_i

Gram matrix
centered about X_i

$$\hat{w}^{(i)} = \underset{w, \lambda}{\operatorname{argmin}} \quad w G^{(i)} w^T + \lambda_i (1 - w \mathbb{1})$$

$$\frac{\partial}{\partial w} : \quad 2 G^{(i)} w^T - \lambda_i \mathbb{1} \stackrel{!}{=} 0 \quad u = \frac{2 w^T}{\lambda}$$

$$G^{(i)} u = \mathbb{1} \quad \text{linear system, solve with standard methods}$$

$$\hat{u}^{(i)} = \left(G^{(i)} \right)^{-1} \cdot \mathbb{1}$$

$$\boxed{\hat{w}^{(i)} = \frac{(\hat{u}^{(i)})^T}{(\hat{u}^{(i)})^T \cdot \mathbb{1}}} \quad \text{repeat for every } i$$

embedding:

$$\hat{z} = \underset{z}{\operatorname{argmin}} \sum_i \| z_i - \hat{w}^{(i)} \tilde{z}^{(i)} \|^2$$

$\tilde{z}^{(i)}$: only neighbors of i

rewrite in matrix notation: $\tilde{w}^{(i)}$: $\hat{w}^{(i)}$ filled with zeros for non-neighbors
 \tilde{w} : matrix of $\tilde{w}^{(i)}$

$$\hat{z} = \underset{z}{\operatorname{argmin}} \quad \| z - \tilde{w} z \|^2$$

$$= \underset{z}{\operatorname{argmin}} \quad \| z^T (\mathbb{I} - \tilde{w})^T (\mathbb{I} - \tilde{w}) z \|^2$$

$$= \underset{z}{\operatorname{argmin}} \quad \| z^T M z \|^2$$

$$M = (\mathbb{I} - \tilde{w})^T (\mathbb{I} - \tilde{w})$$

$$= U \Lambda U^T \quad \text{eigen-decomp.}$$

assume that eigenvalues are sorted $\Rightarrow u_D$ is the eigenvector for λ_D ,
smallest eigenvalue

$\lambda_D = 0$ because M has rank $D-1$ ($M \cdot \mathbf{1} = 0 \cdot \mathbf{1}$)

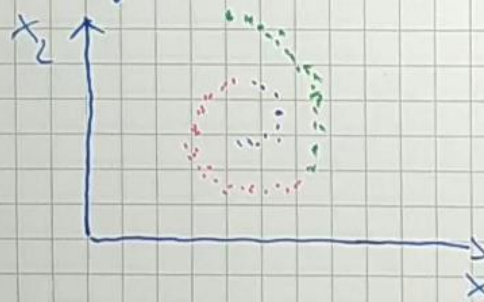
$\lambda_{D-1} \neq 0 \Rightarrow$ this minimizes the objective

when $\dim(Z) = D' \Rightarrow$ use the columns $u_{D-D'+1}, \dots, u_{D-1}$ as new features

$$z_{ij} = \frac{u_{i, D-D'+1+j}}{u_{i, D-D'+1+j}}$$

- embed new data: x_{new} :
 - find the k nearest neighbors of x_{new} in TS
 - find w_{new} such that $w_{\text{new}} = \arg\min_w \|x_{\text{new}} - w \tilde{x}^{(new)}\|^2$
 \uparrow
 neighbors
 - define $z^{(new)} = w_{\text{new}} \cdot \tilde{z}^{(new)}$
 \uparrow
 neighbors of x_{new} in z -space

- example



LLE unrolls the spiral
(k -PCA cannot do this!)

- in practice: t -SNE is most popular $\hat{=}$ improvement of LLE

new instance x_{new} : row vector of size D

w_{new} : row vector of size k

z_{new} : row vector of size D' , $\tilde{z}^{(new)} \in \mathbb{R}^{k \times D'}$

$\tilde{x}^{(new)} \in \mathbb{R}^{k \times D}$
 \uparrow
 k rows for k nearest neighbors.

$\tilde{x}_i^{(new)} \in TS$ s.t. $d(x_{\text{new}}, \tilde{x}_i^{(new)})$ is

small
 \leftarrow new features of k nearest neighbors from TS