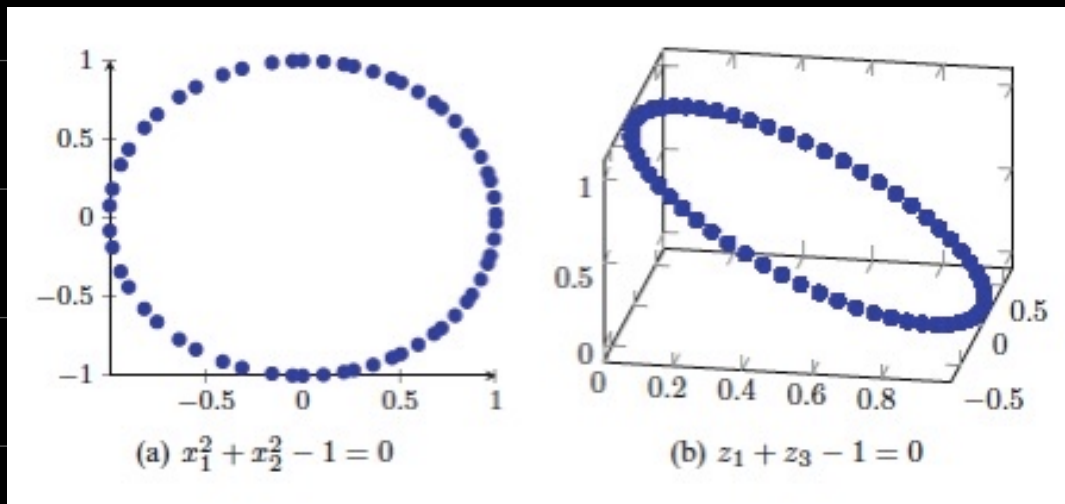


## 3.4 Nonlinear / Kernel PCA

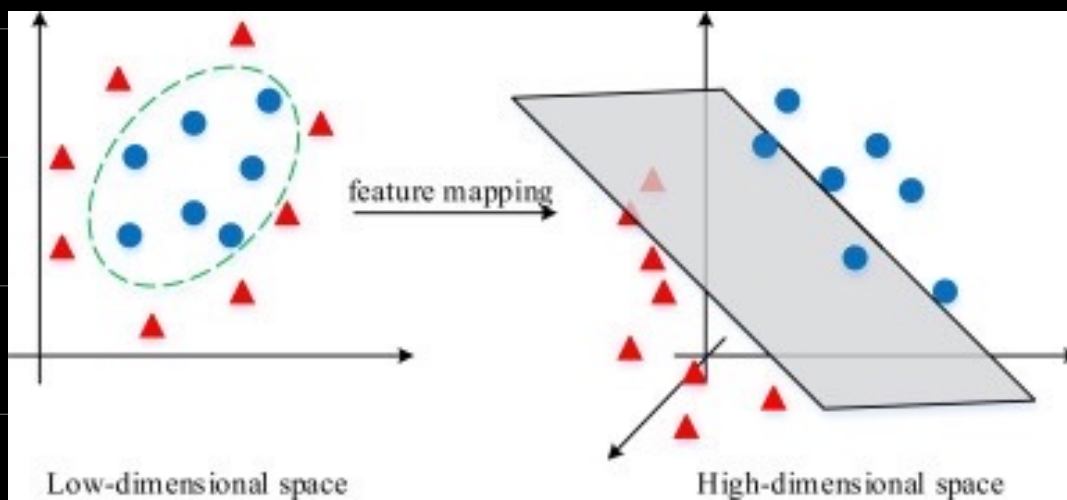
### 3.4.1 Motivation



$$\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

Observation: (1) nonlinear change of features can turn nonlinear eqns into linear eqns.



Observation: (2) in some cases, increasing the number of features is helpful.

Nonlinear/Kernel PCA: a method to embed sample pts into a high dimensional space before applying PCA such that the structure of the sample becomes linear.

Setting: Given sample pts  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$  consider a nonlinear embedding map  $\phi$ , known as the feature map,

$$\phi: \mathbb{R}^p \rightarrow \mathbb{R}^D \quad p \ll D$$

$$\begin{array}{ccccccc} \checkmark & x^{(i)} = & \underbrace{(x_1^{(i)}, \dots, x_p^{(i)})^T}_{\text{old features}} & \mapsto & \underbrace{\phi(x^{(i)}) = (\phi_1(x^{(i)}), \dots, \phi_D(x^{(i)}))^T}_{\text{new features}} \\ \text{ith sample pt} & & & \text{image of} & & & \\ & & & \text{the ith sample pt} & & & \end{array}$$

In the higher dim space  $\mathbb{R}^D$ , the sample pts become  $\phi(x^{(1)}), \dots, \phi(x^{(n)}) \in \mathbb{R}^D$ .

and the sample matrix becomes

$$\Phi \stackrel{\text{def}}{=} \begin{pmatrix} | & & | \\ \phi(x^{(1)}) & \dots & \phi(x^{(n)}) \\ | & & | \end{pmatrix}_{D \times n}$$

In nonlinear / kernel PCA, we do:

Step 1: center  $\phi(x^{(1)}), \dots, \phi(x^{(n)})$ .

If  $\mu \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \phi(x^{(i)}) \neq 0$ , replace them by  $\phi(x^{(1)}) - \mu, \dots, \phi(x^{(n)}) - \mu$ .

This corresponds to replace  $\Phi$  by  $\Phi H$  where  $H \stackrel{\text{def}}{=} I - \frac{1}{n} \mathbf{1} \mathbf{1}^T$  is the centering matrix.

Step 2: apply PCA to  $\phi(x^{(1)}), \dots, \phi(x^{(n)})$ .

i.e., find the eigenvectors associated to the largest eigenvalues of  $\Phi \Phi^T$

However,  $\underbrace{\Phi}_{D \times n} \underbrace{\Phi^T}_{n \times D} \in \mathbb{R}^{D \times D}$ , it can be computationally expensive to get its EVD if  $D$  is large.

### 3.4.2. Math Prep.

Let  $A \in \mathbb{R}^{m \times n}$ .

Lemma: If  $v$  is an eigenvector of  $A^T A$  associated to  $\lambda > 0$ , then  
$$u \stackrel{\text{def}}{=} Av$$

is an eigenvector of  $AA^T$  associated to the same  $\lambda$ . Moreover,  $\|u\| = \sqrt{\lambda} \|v\|$ .

Proof: We have  $A^T A v = \lambda v$ , then

$$AA^T u \stackrel{\text{def of } u}{=} A(A^T A v) \stackrel{v \text{ is eigenvector}}{=} A(\lambda v) = \lambda Av \stackrel{\text{def of } u}{=} \lambda u$$

which means  $u$  is an eigenvector of  $AA^T$  associated to  $\lambda$ . Moreover,

$$\overset{\text{def of } u}{\|u\|^2} = \|Av\|^2 = (Av)^T(Av) = v^T(A^TAv) \overset{v \text{ is an eigenvector}}{=} v^T(\lambda v) \\ = \lambda \|v\|^2.$$

### 3.4.3. Theory.

If  $n \ll D$ , the lemma says that the eigenvectors of  $\Phi\Phi^T \in \mathbb{R}^{D \times D}$  can be obtained from the eigenvectors of  $\Phi^T\Phi \in \mathbb{R}^{n \times n}$ .

We proceed in this way to reduce the computational load.

$$\text{Let } (\Phi^T\Phi)_{n \times n} = U \Lambda U^T = \begin{pmatrix} | & & | \\ u^{(1)} & \dots & u^{(n)} \\ | & & | \end{pmatrix} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \begin{pmatrix} \overbrace{u^{(1)}} \\ \vdots \\ \underbrace{u^{(n)}} \end{pmatrix}$$

where  $\lambda_1 \geq \dots \geq \lambda_n \geq 0$  and  $u^{(i)}$  is a unit eigenvector of  $\Phi^T\Phi$  associated to  $\lambda_i$ .

By the lemma,  $\Phi u^{(i)}$  is an eigenvector of  $\Phi\Phi^T$  associated to  $\lambda_i$ , and

$$\|\Phi u^{(i)}\| = \sqrt{\lambda_i} \|u^{(i)}\| = \sqrt{\lambda_i}.$$

Thus  $v^{(i)} \stackrel{\text{def}}{=} \frac{1}{\sqrt{\lambda_i}} \Phi u^{(i)}$  is a unit eigenvector of  $\Phi\Phi^T$  associated to  $\lambda_i$ .