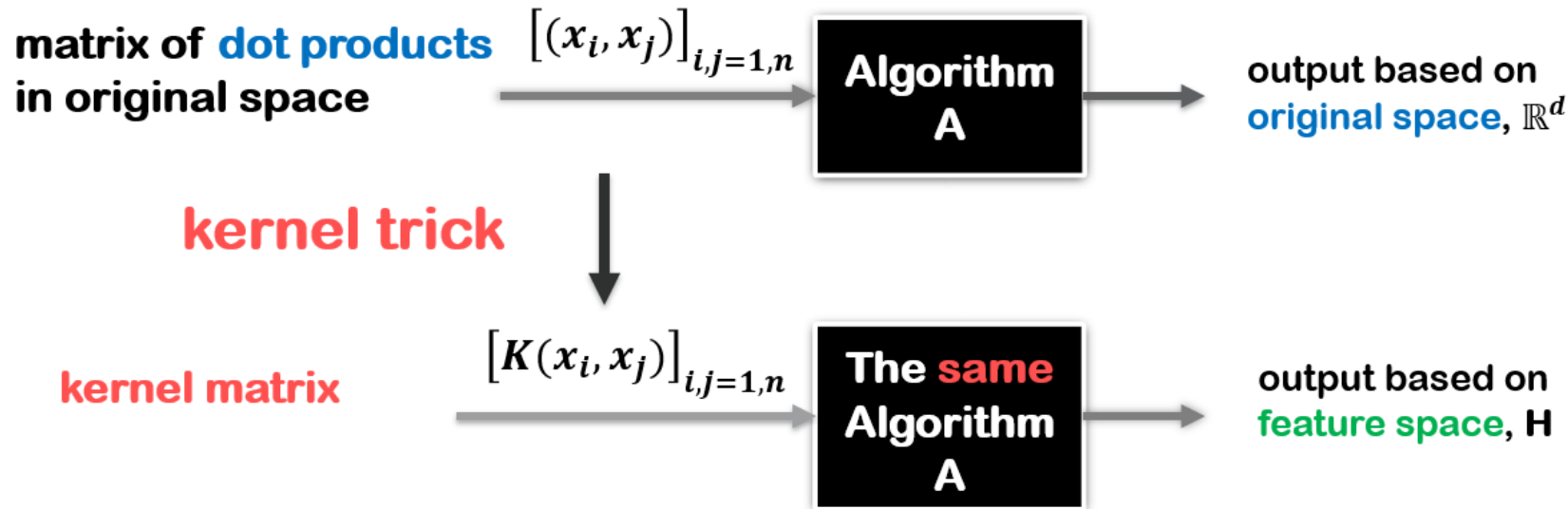


# Kernels and Kernel Tricks: Non-linearity



- **Kernel**
- **Kernel Matrix**
- **Kernel Trick**
- **Kernelization**

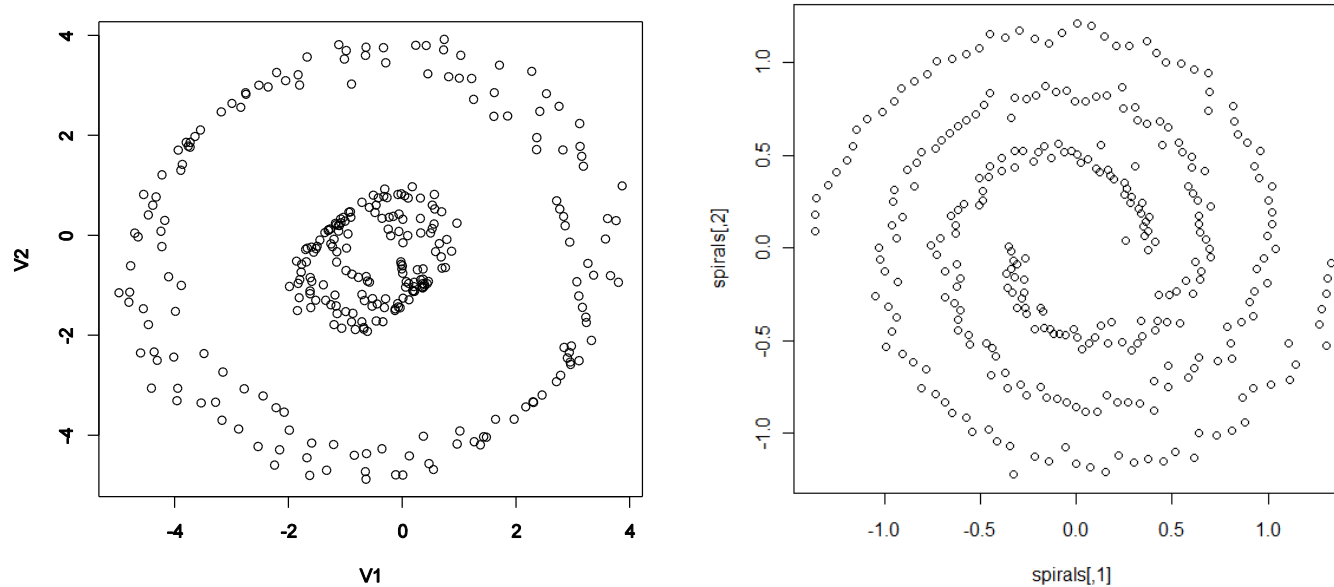
**Prof. Nagiza F. Samatova**

samatova@csc.ncsu.edu

Department of Computer Science  
North Carolina State University

# Motivation: PCA-based Dimension Reduction

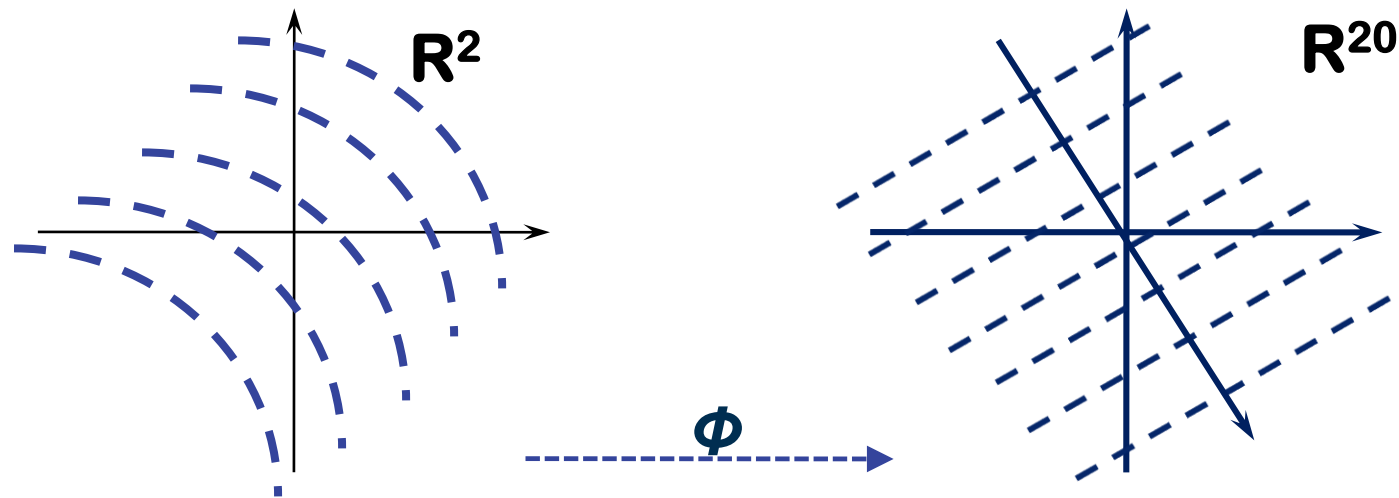
- Standard PCA does not work for all datasets



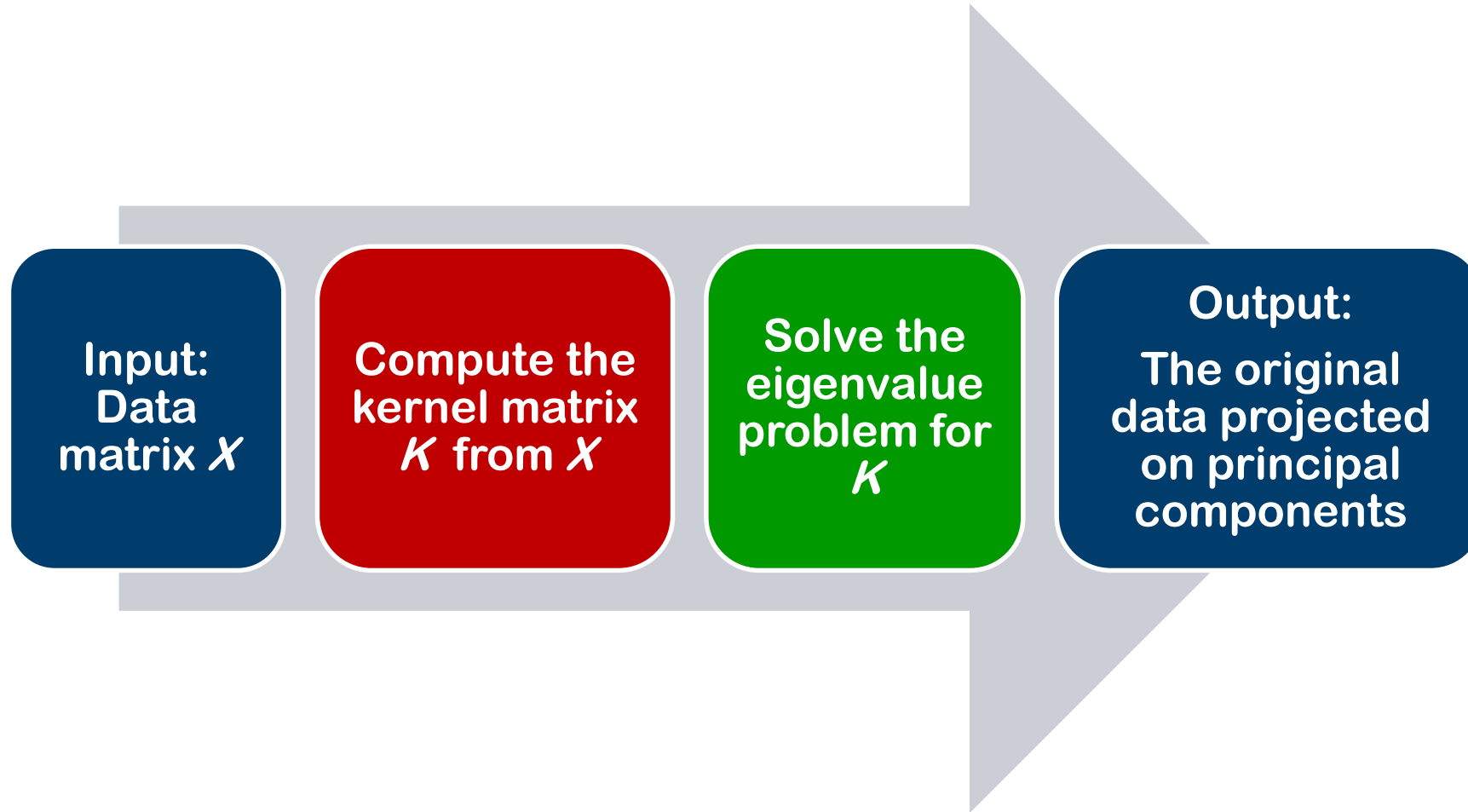
- Absence of a hyper-plane or principal components to separate data clearly

# Solution: Input space $\rightarrow$ Feature space

- Map data objects in **non-linear input space** to a **linear feature space** (Hilbert)
  - Potentially increases dimensionality of the data



# Kernel-PCA



# What was the Trick?

- Kernel-PCA performs PCA in a feature space **without explicitly mapping** the input space to a feature space!

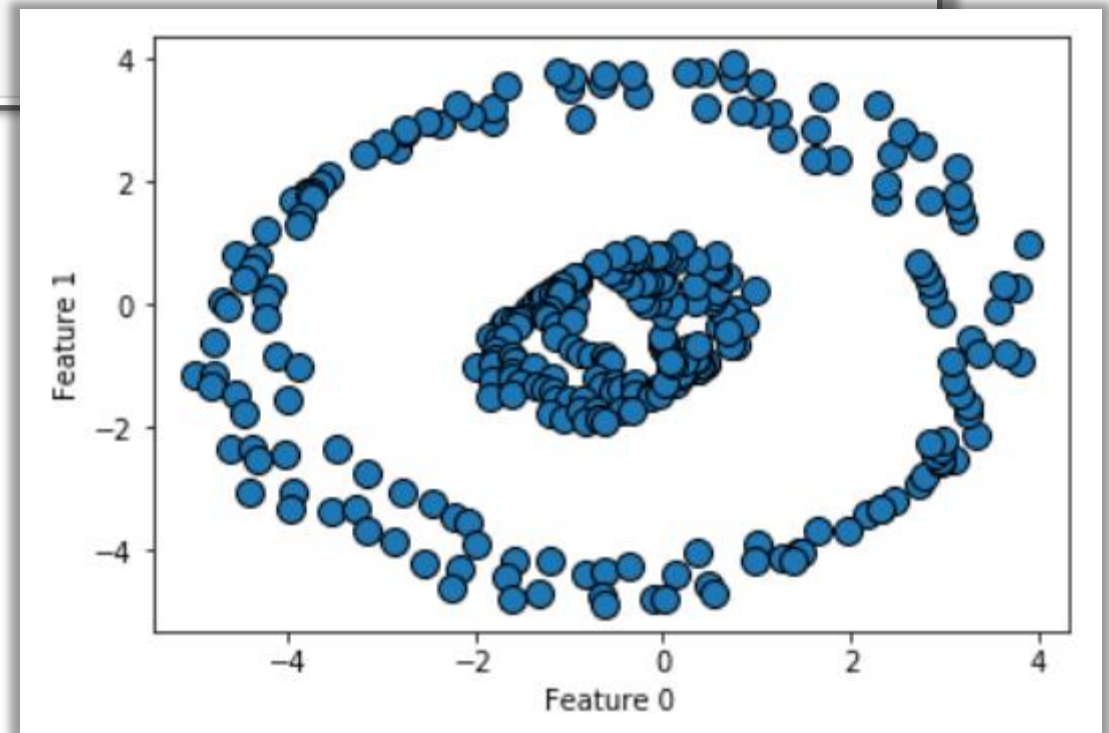
DR technique without a kernel	DR using a <b>kernel</b>
Requires an explicit mapping	Does not require an explicit mapping
Maps each data object to feature space	Does not explicitly map each data object
Requires a covariance matrix in feature space	Requires a <b>kernel matrix</b>
Solves the eigenvalue problem for the covariance matrix	Solves the eigenvalue problem on the <b>kernel matrix</b>

# Code: Load and Plot the Non-Linear Data

```
1 import mglearn
2 import pandas as pd
3 circles = pd.read_csv("../data_raw/data_prep_kernels_circles.csv")
4 mglearn.discrete_scatter(circles.iloc[:,0], circles.iloc[:,1])
5 plt.xlabel("Feature 0")
6 plt.ylabel("Feature 1")
7 plt.show()
```

## Linear dimension reduction:

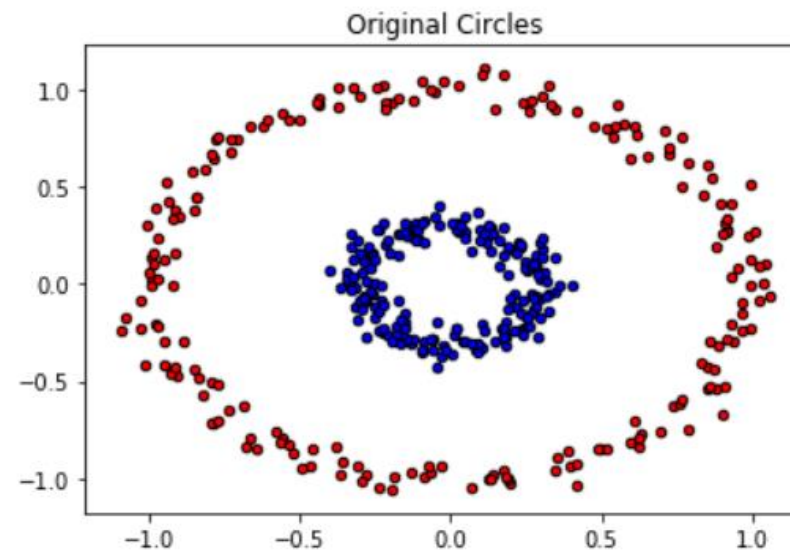
- any projection onto Feature-0 axis or Feature-1 will fail to separate the points on two different circles



# Kernel PCA

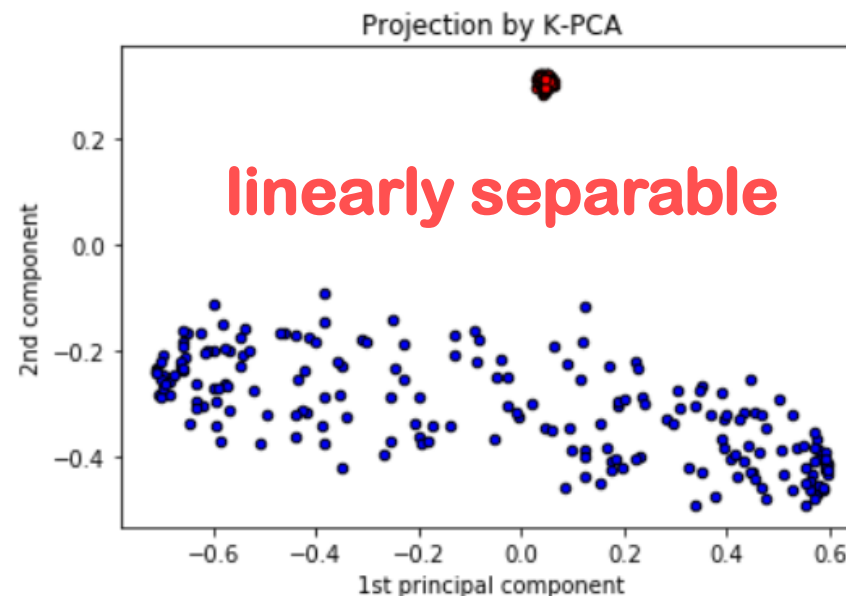
```
from sklearn.datasets import make_circles
np.random.seed(0)
X, y = make_circles(n_samples=400, factor=.3, noise=.05)
```

```
plt.figure()
plt.title("Original Circles")
reds = y == 0
blues = y == 1
plt.scatter(X[reds, 0], X[reds, 1], c="red",
            s=20, edgecolor='k')
plt.scatter(X[blues, 0], X[blues, 1], c="blue",
            s=20, edgecolor='k')
plt.show()
```



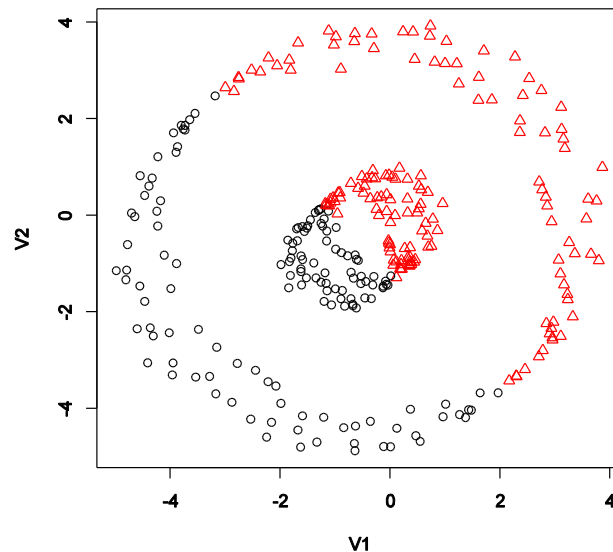
```
from sklearn.decomposition import KernelPCA
kpca = KernelPCA(kernel="rbf", fit_inverse_transform=True, gamma=10)
X_kpca = kpca.fit_transform(X)
```

```
plt.scatter(X_kpca[reds, 0], X_kpca[reds, 1], c="red",
            s=20, edgecolor='k')
plt.scatter(X_kpca[blues, 0], X_kpca[blues, 1], c="blue",
            s=20, edgecolor='k')
plt.title("Projection by K-PCA")
plt.xlabel("1st principal component")
plt.ylabel("2nd component")
plt.show()
```



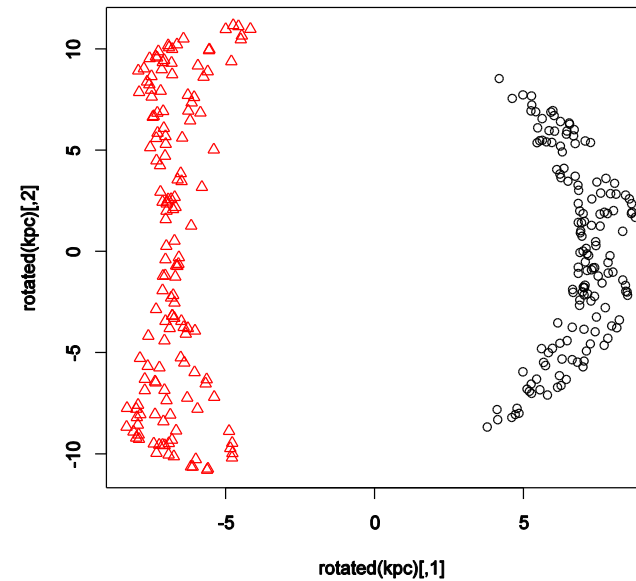
# Example: Kernel-PCA assisting Clustering

k-means clustering **before**  
kernel-PCA



**Incorrect Clusters**

k-means clustering **after**  
kernel-PCA



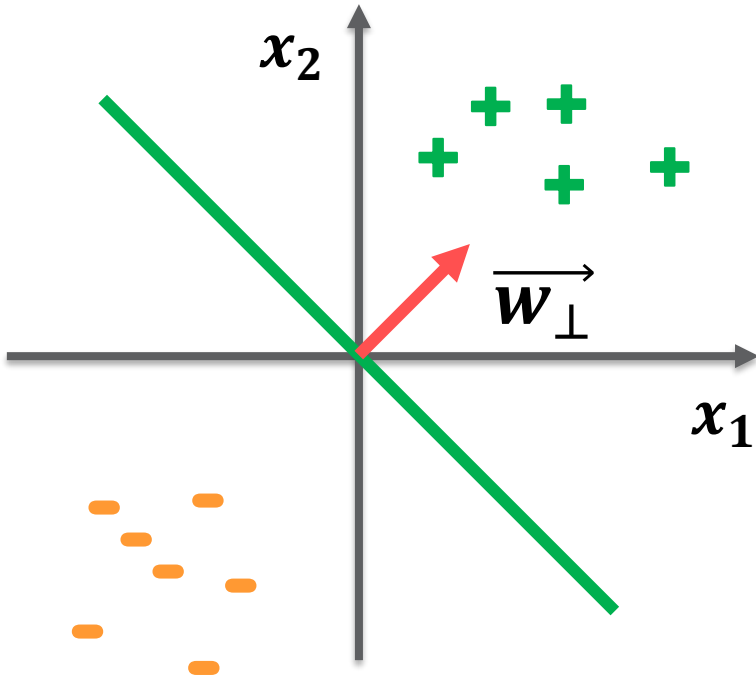
**Correct Clusters**



# Kernels

## **UNDERLYING MATH CONCEPTS**

# Preliminaries: **Linear** Decision Boundary



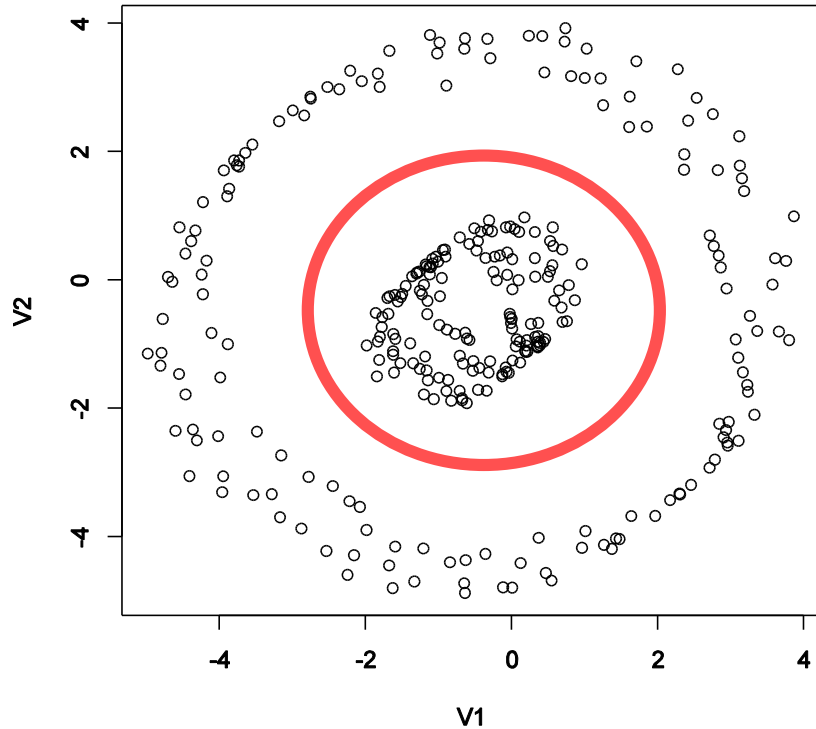
1.  $\langle \vec{u} \cdot \vec{v} \rangle = \vec{u} \cdot \vec{v} = (\vec{u}, \vec{v})$  is a **similarity measure** between points/vectors (**scalar product, inner product, dot product**)

2.  $\vec{w}_\perp \cdot \vec{x} + b = 0$  – equation of a hyper-plane w/ the normal (perpendicular) vector  $\vec{w}_\perp$

3.  $\vec{w}_\perp \cdot \vec{x} + b \geq 0$  – **positive** hyper-space  
 $\vec{w}_\perp \cdot \vec{x} + b < 0$  – **negative** hyper-space

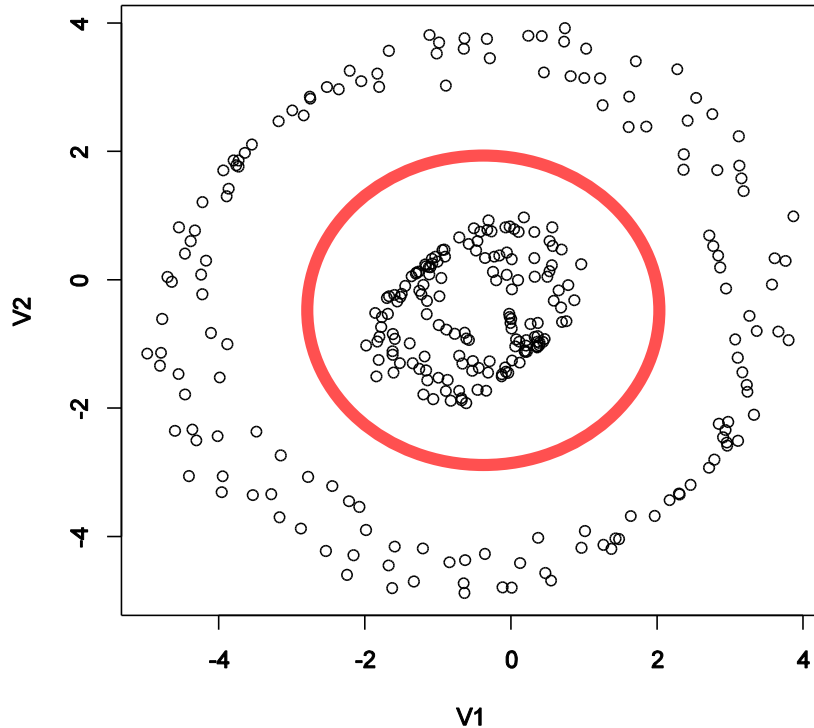
Thus, if a hyper-plane separates (**linearly**) different classes of points, then  $(\vec{w}_\perp, b)$  is all we need to know to decide the class for an unknown object  $\hat{x}$ :  $\text{sign}(\vec{w}_\perp \cdot \vec{x} + b)$ ,  $\text{sign} = (+/-)$

# Non-linear Decision Boundary



Solving **non-linear** problems is usually much harder than solving linear optimization problems.

# Explicit Mapping to Deal with Non-linearity



Imagine that we could transform such non-linearly separable data into some higher-dimensional space, where **transformed data becomes linearly separated**:

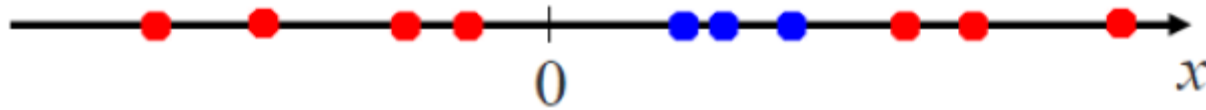
- reduce a **non-linear problem** in the original space to the **linear problem** in a higher-dimensional, transformed space.

## Why transform?

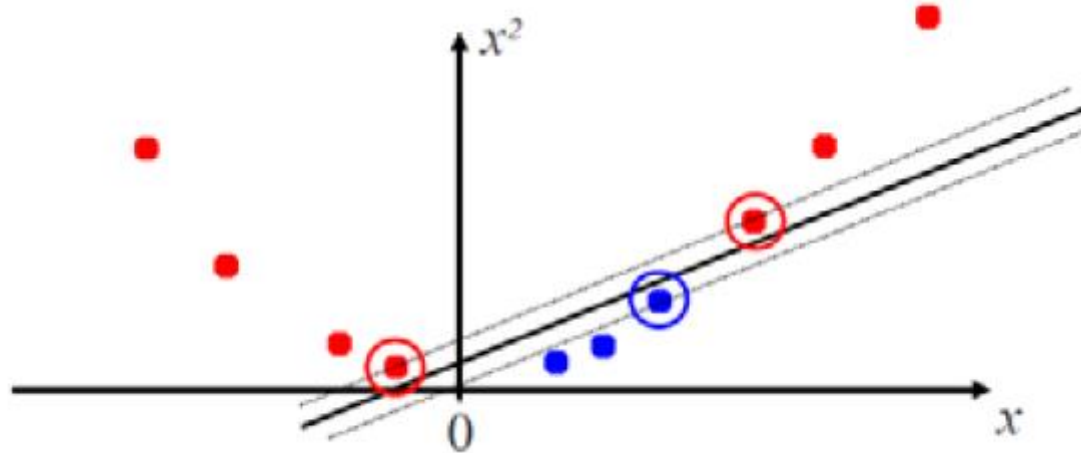
- Linear operation in the feature space is equivalent to non-linear operation in input/original space
- Classification can become easier with a proper transformation

# Linear Separability Example

- Is this data linearly separable?



- How about a **quadratic** mapping:  $\Phi(x^2)$



# Ex #1: Explicit Mapping

Example:  $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

$$\mathbf{p} = (x_1, x_2) \in \mathbb{R}^2 \rightarrow \mathbf{p}_{new} = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \in \mathbb{R}^3$$

$$\mathbf{u} = (1, -2)$$

$$\mathbf{v} = (2, 3)$$

inner / dot / scalar product

$$\begin{aligned} (\mathbf{u}, \mathbf{v}) &= \\ &= 1 * 2 + (-2) * 3 \\ &= -4 \end{aligned}$$

$$\mathbf{u}_{new} = (1, -2\sqrt{2}, 4) \in \mathbb{R}^3$$

$$\mathbf{v}_{new} = (4, 6\sqrt{2}, 9) \in \mathbb{R}^3$$

inner / dot / scalar product

$$\begin{aligned} (\mathbf{u}_{new}, \mathbf{v}_{new}) &= \\ &= 1 * 4 + (-12) * 2 + 4 * 9 \\ &= 4 - 24 + 36 = 16 \\ &= (\mathbf{u}, \mathbf{v})^2 \end{aligned}$$

# Example: Explicit Mapping

Example:  $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

$$\mathbf{p} = (x_1, x_2) \in \mathbb{R}^2 \rightarrow \mathbf{p}_{new} = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \in \mathbb{R}^3$$

New coordinates are **explicitly** defined through original coordinates:

$$\begin{aligned}\overrightarrow{\Phi(\mathbf{u})} &= \mathbf{u}_{new} = (u_1^2, \sqrt{2}u_1u_2, u_2^2) \\ \overrightarrow{\Phi(\mathbf{v})} &= \mathbf{v}_{new} = (v_1^2, \sqrt{2}v_1v_2, v_2^2)\end{aligned}$$

Since  $(\vec{u}, \vec{v})$ , the **inner product**, is a **similarity** between points in original space:

- what happens with the similarity between newly transformed points:

$$\begin{aligned}(\overrightarrow{\Phi(\mathbf{u})}, \overrightarrow{\Phi(\mathbf{v})}) &= (\mathbf{u}_{new}, \mathbf{v}_{new}) = \\ u_1^2v_1^2 + 2u_1u_2v_1v_2 + v_1^2v_2^2 &= \\ (u_1v_1 + u_2v_2)^2 &= \langle \vec{u} \cdot \vec{v} \rangle^2\end{aligned}$$

# Kernel Function = F (scalar products)

- **Kernel function in terms of:**

- the original inner/scalar/dot products
- without explicit knowledge of new coordinates!

E.g.:  $K(u, v) = (u, v)^2$

- **No need to know the coordinates of vectors  $u$  and  $v$**

- only their similarity in terms of inner product=scalar product = dot product.

- **It is very handy when**

- the objects are strings, graphs, images, texts
- for which it is easier to define similarity metric
- rather than coordinates of their vector representation



# Implicit Mapping (Kernel) vs. Explicit Mapping

Explicit mapping (not desirable):

$\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^m = H(m > d)$  – feature space, Hilbert space (generalization of a Euclidian space),  $m = \infty$  possible such that:

$$\forall \vec{u}, \vec{v} \in \mathbb{R}^d: (\vec{u}, \vec{v}) \approx (\Phi(\vec{u}), \Phi(\vec{v}))$$

inner products/similarity in transformed space are similar to inner products in original space

Implicit representation:

$$K(u, v) = (\Phi(u), \Phi(v)) = F((u, v))$$

is the inner product in a Hilbert space,  $H$ , i.e.,

no need to know coordinates:  $\Phi(u) \in \mathbb{R}^m = H$

# Examples: Kernel Functions

$$K(u, v) = (u, v)^2 - \text{vanilladot kernel}$$

$$K(u, v) = (u, v)^2 \text{ or}$$

$$K(u, v) = (\text{scale} * (u, v) + \text{offset})^{\text{degree}} - \text{polynomial kernel}$$

$$K(u, v) = e^{-\frac{|u-v|^2}{\sigma}} - \text{Radial Basis Function (RBF)}$$

$(\mathbb{R}^m = H, m = \infty \text{ infinite Hilbert space})$

Euclidean distance:

Inner/dot/scalar product

$$|u - v|^2 = (u - v, u - v) = (u, u) - 2(u, v) + (v, v)$$

# Formal Definition of a Kernel

$K(u, v)$  is a **kernel** if

(1) **Symmetric**:  $K(u, v) = K(v, u)$  (\*)

(2) **Positive semi-definite**:

$\forall c_1, c_2, \dots, c_n \in \mathbb{R}$  and a set of points

$\forall u_1, u_2, \dots, u_n \in \mathbb{R}^d$

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j K(u_i, u_j) \geq 0 \quad (**)$$

- valid **similarity measure**

In terms of matrix notation (\*\*) can be written as:

$C^T \cdot M_K \cdot C \geq 0$ , where

$$C = \begin{bmatrix} c_1 \\ c_2 \\ \dots \\ c_n \end{bmatrix} \text{ and } \mathbf{M}_K = \begin{bmatrix} K(u_1, u_1) & K(u_1, u_2) & \dots & K(u_1, u_n) \\ K(u_2, u_1) & K(u_2, u_2) & \dots & K(u_2, u_n) \\ \dots & \dots & \dots & \dots \\ K(u_n, u_1) & K(u_n, u_2) & \dots & K(u_n, u_n) \end{bmatrix} \quad \text{kernel matrix}$$

# Euclidean Distance as Dot Product

Let  $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{n \times d}$  be a set of points in  $d$ -dimensional space.

It is a **vector** data. Let  $(x_i, x_j)$  be the dot product between two vectors.

Note: that **dot (inner/scalar) product is a similarity measure**

Distance (**Euclidean**) between  $x_i$  and  $x_j$  can be calculated in terms of **dot/inner products**:

$$\begin{aligned} d^2(x_i, x_j) &= \left\| x_i - x_j \right\|^2 = (x_i - x_j, x_i - x_j) \\ &= (x_i, x_i) - 2(x_i, x_j) + (x_j, x_j) \end{aligned}$$

# Kernel Trick and Kernelization

Distance (Euclidian) between  $x_i$  and  $x_j$ :

$$d^2(x_i, x_j) = (x_i, x_i) - 2(x_i, x_j) + (x_j, x_j)$$

What happens if we replace dot product by a kernel evaluation?

$x_i \rightarrow \Phi(x_i)$  and  $x_j \rightarrow \Phi(x_j)$  in Hilbert space (Generalized Euclidean space)

$$d^2(\Phi(x_i), \Phi(x_j)) = K(x_i, x_i) - 2K(x_i, x_j) + K(x_j, x_j)$$

$$\begin{aligned} d^2(\Phi(x_i), \Phi(x_j)) &= \|\Phi(x_i) - \Phi(x_j)\|^2 = \\ &= (\Phi(x_i) - \Phi(x_j), \Phi(x_i) - \Phi(x_j)) = \\ &= (\Phi(x_i), \Phi(x_i)) - 2(\Phi(x_i), \Phi(x_j)) + (\Phi(x_j), \Phi(x_j)) = \\ &= K(x_i, x_i) - 2K(x_i, x_j) + K(x_j, x_j) \end{aligned}$$

$$\begin{aligned} d^2(x_i, x_j) &= \|x_i - x_j\|^2 = (x_i - x_j, x_i - x_j) \\ &= (x_i, x_i) - 2(x_i, x_j) + (x_j, x_j) \end{aligned}$$

## Ex #2: Euclidean Distance as Dot Products

Example:  $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

$$\mathbf{p} = (x_1, x_2) \in \mathbb{R}^2 \rightarrow \mathbf{p}_{new} = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \in \mathbb{R}^3$$

$$\mathbf{u} = (1, -2)$$

$$\mathbf{v} = (2, 3)$$

$$d^2(\mathbf{u}, \mathbf{v}) =$$

$$\mathbf{u}_{new} = (1, -2\sqrt{2}, 4) \in \mathbb{R}^3$$

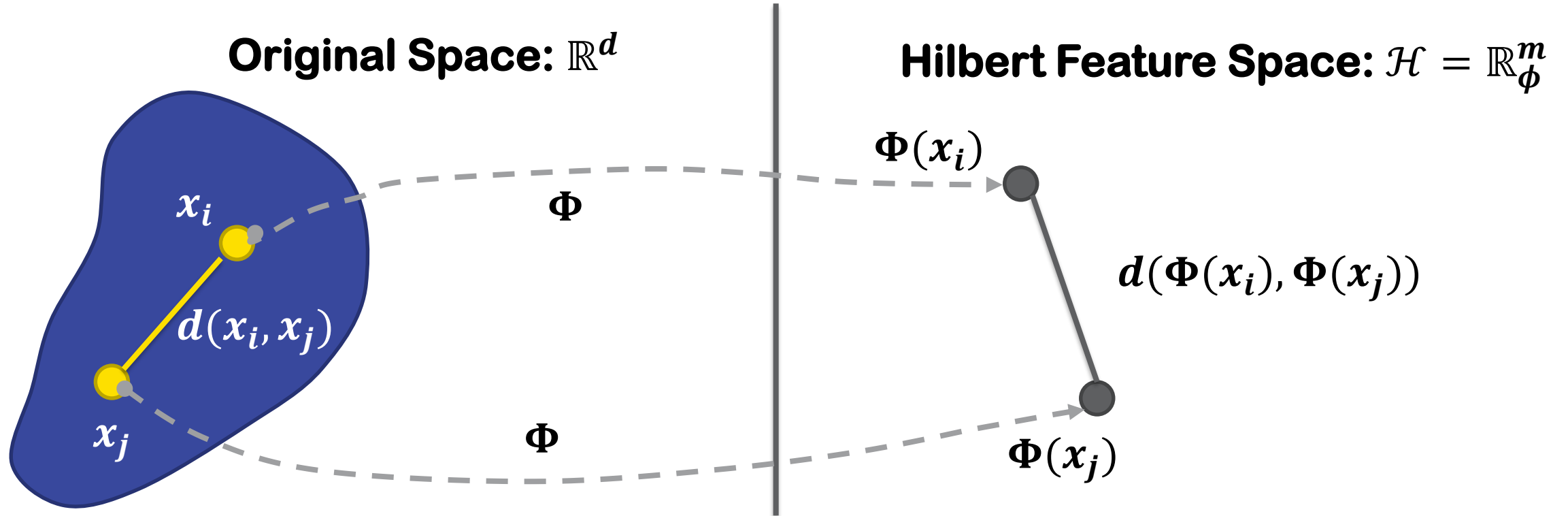
$$\mathbf{v}_{new} = (4, 6\sqrt{2}, 9) \in \mathbb{R}^3$$

$$d^2(\mathbf{u}_{new}, \mathbf{v}_{new}) =$$

# Euclidean Distance in H

$$d^2 \left( \Phi(x_i), \Phi(x_j) \right) = K(x_i, x_i) - 2K(x_i, x_j) + K(x_j, x_j)$$

b/s by def. of a kernel  $K(u, v) = (\Phi(u), \Phi(v))$

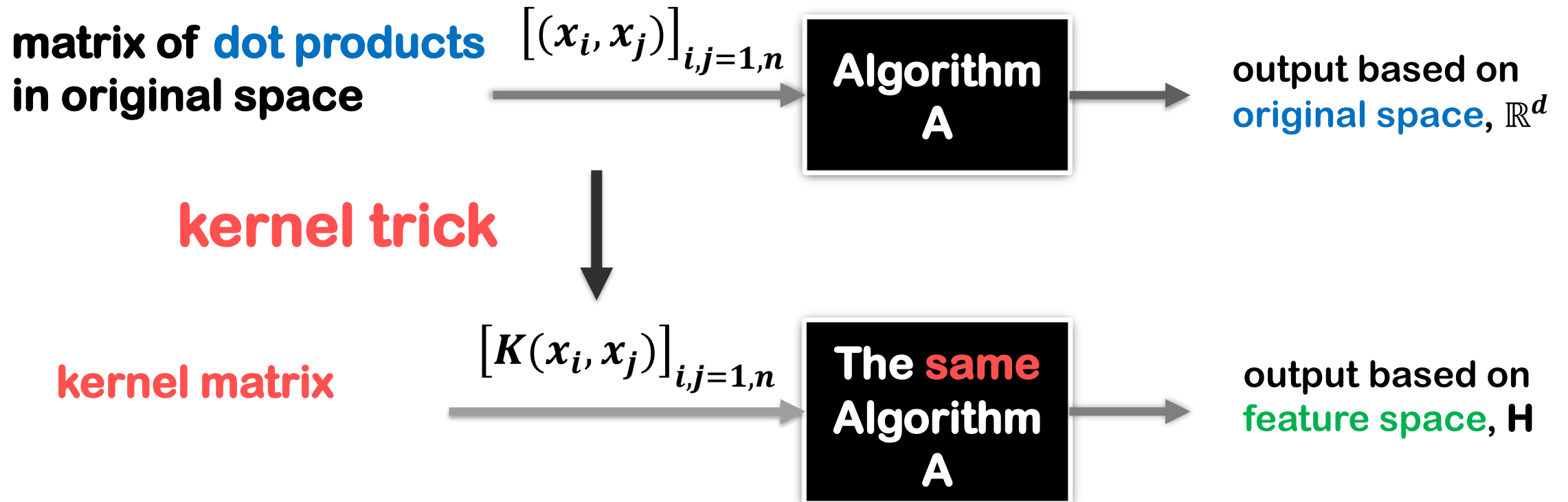


# Kernelization

If  $K(u, v) = F((u, v))$  (e.g.,  $K(u, v) = (u, v)^2$ ),  
then we can compute distances (similarities) in the transformed (Hilbert) space

(a) only using **dot products** of original data

(b) **without any explicit knowledge** of vector coordinates in the feature space  $H$

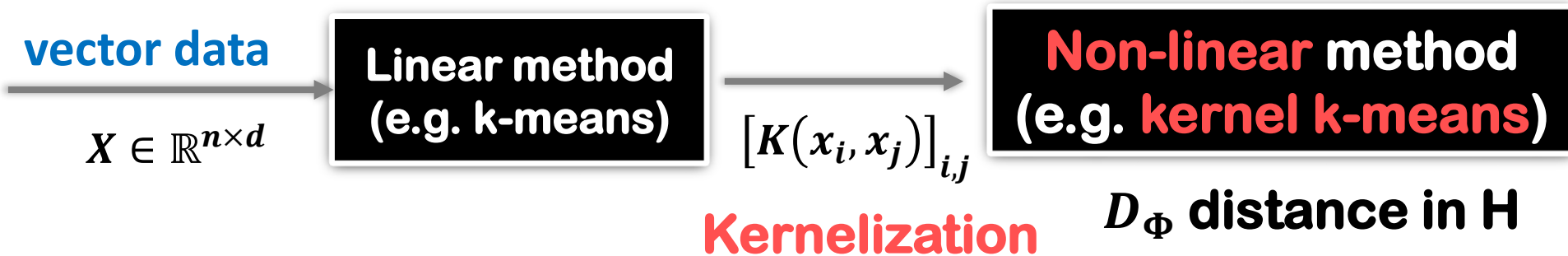
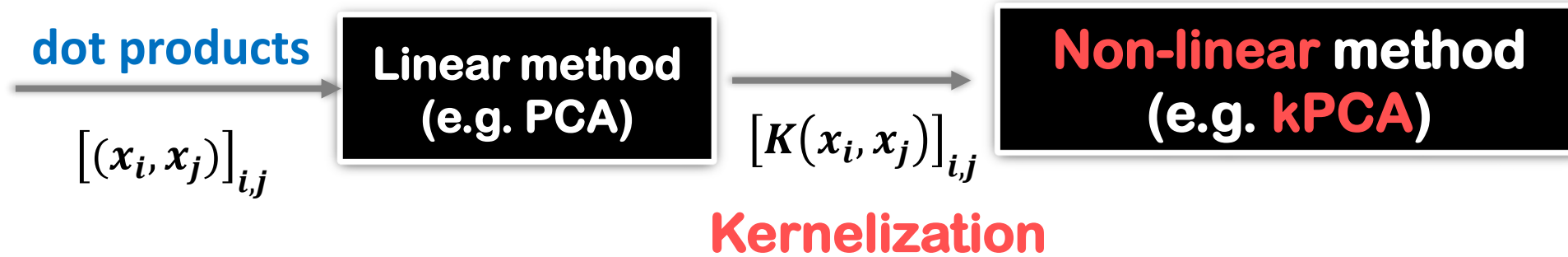




# Kernel Trick and Kernelization

- Thus, the **kernel trick** will
  - replace **dot products** over vector data in original space  $\mathbb{R}^d$
  - with their kernel evaluations ( $K(u, v) = F((u, v))$ )
    - expressed as symmetric, semi-definite  $F()$  over  $(u, v)$
- and will apply the **same** algorithm to the **kernel matrix**
- Basically, the **kernel trick** will
  - generalize **linear methods** (e.g., PCA) to **non-linear** methods (e.g., kernel PCA)
  - by simply replacing classic **dot products** by a **kernel**
- This transformation is called **kernelization**

# Examples: Kernel Tricks



# Kernels and Kernel Tricks

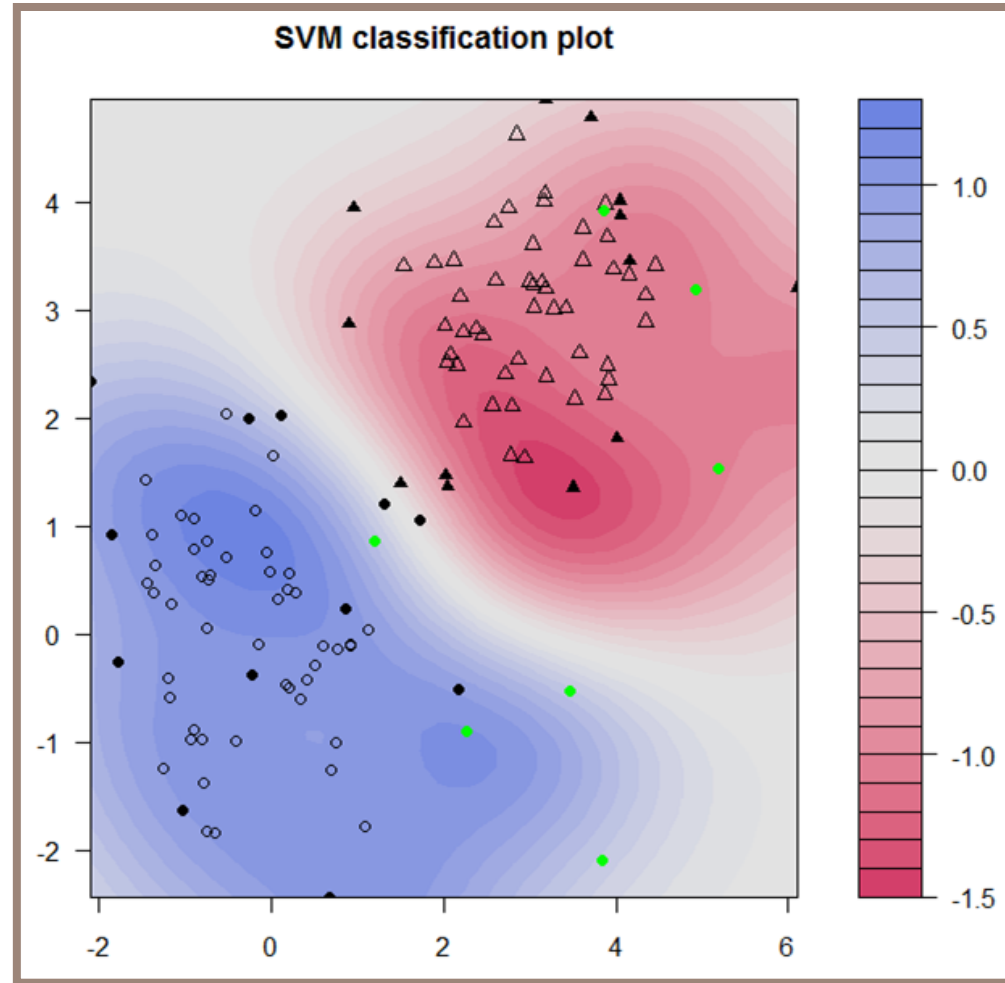
## **CLASSIFICATION**

# Classification with Kernel SVM

- **SVM Optimization problem:**
  - The data points only appear in the **inner products**
- **As long as we can calculate the inner product in the feature space:**
  - we do not need any **explicit mapping**
- **Common **geometric operations**:**
  - angles (cosine angle)
  - distances
- **They can be expressed by **inner products****
- **Define the kernel function as the function of inner products**
  - this give the **implicit mapping**

$$\begin{aligned} \max_{\alpha} \quad & \mathcal{J}(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned}$$

# Example: Linear Separable SVM

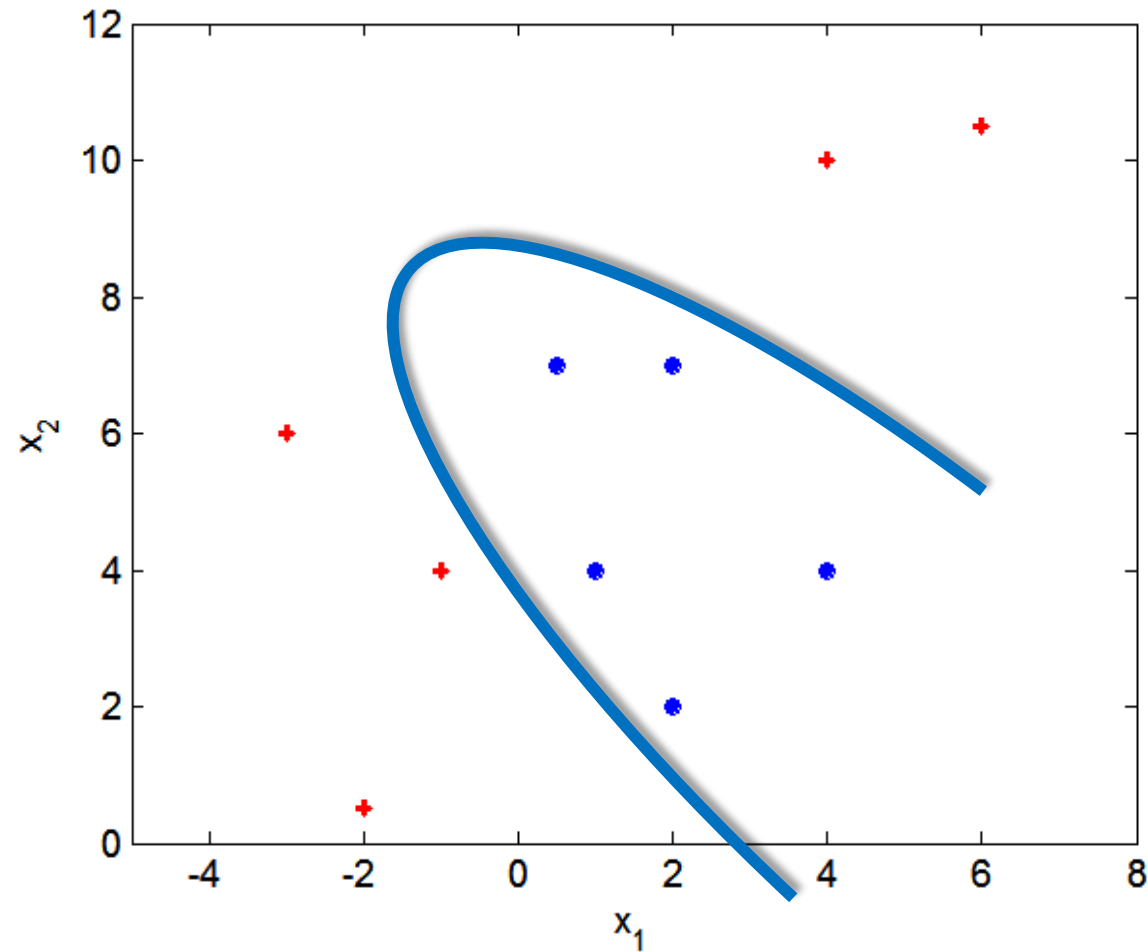


# Kernels

## **NON-LINEAR PROBLEMS**

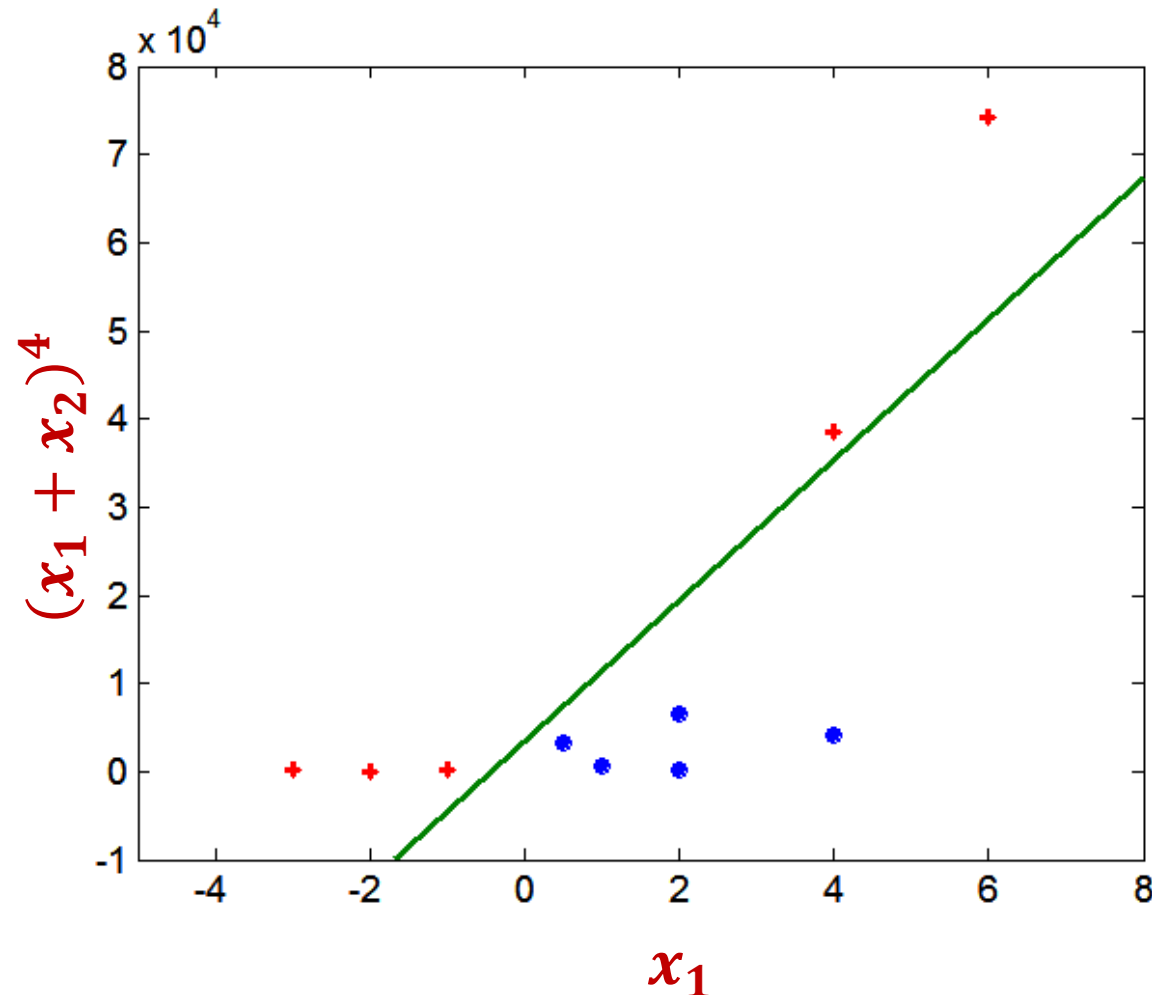
# Nonlinear Support Vector Machines

- What if the **decision boundary** is **not-linear**?



# Non-linear Support Vector Machines

Transformed data into higher dimensional space  
using Kernel functions





# Example: Nonlinear SVM

```
1. library ("kernlab")
2. help (ksvm)
3. x = rbind(matrix(rnorm(120),,2),matrix(rnorm(60,mean=3),,2),
      matrix(rnorm(60, mean=-3), , 2))
4. y = matrix(c(rep(1,60), rep(-1,60)))
5. cl = c (rep("purple",60), rep("red",60))
6. plot (x, col=cl, pch=19)

7. # play with the cost of constraints
8. model = ksvm (x, y, type="C-svc", C=0)
9. model; error (model)
10. model = ksvm (x, y, type="C-svc", C=1)
11. model

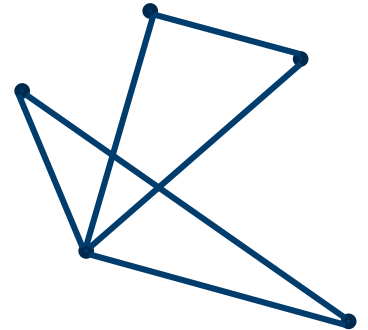
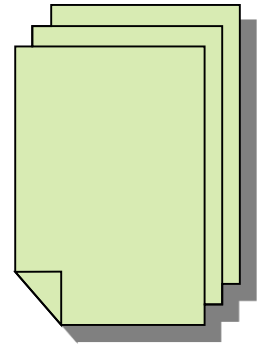
12. # play with kernel functions
13. model = ksvm (x, y, type="C-svc", kernel = "polydot")
14. model
15. model = ksvm (x, y, type="C-svc", kernel = "rbfdot")
16. model
```

# Kernels

## **COMPLEX NON-VECTOR DATA**

# Kernels are ubiquitous

- There exist kernels for a variety of data formats
  - vectors
  - graphs
  - music
  - strings
  - text
  - and so on
- There can be multiple kernels for each data modality:
  - Choosing the right one is an art!



# Summary: Kernels: Strengths & Weaknesses

## ● Pros:

- Generalizes many algorithms that work solely in linear space (non-linear in input space  $\square$  linear in feature space)
- Very simple:
  - does not involve non-linear optimization
  - essentially, linear algebra

## ● Cons:

- Time-consuming for large datasets because of the kernel function computation
- Kernel matrix is of size  $n \times n$  where  $n$  is the number of observations, that becomes a quadratic cost, prohibitive for large values of  $n$