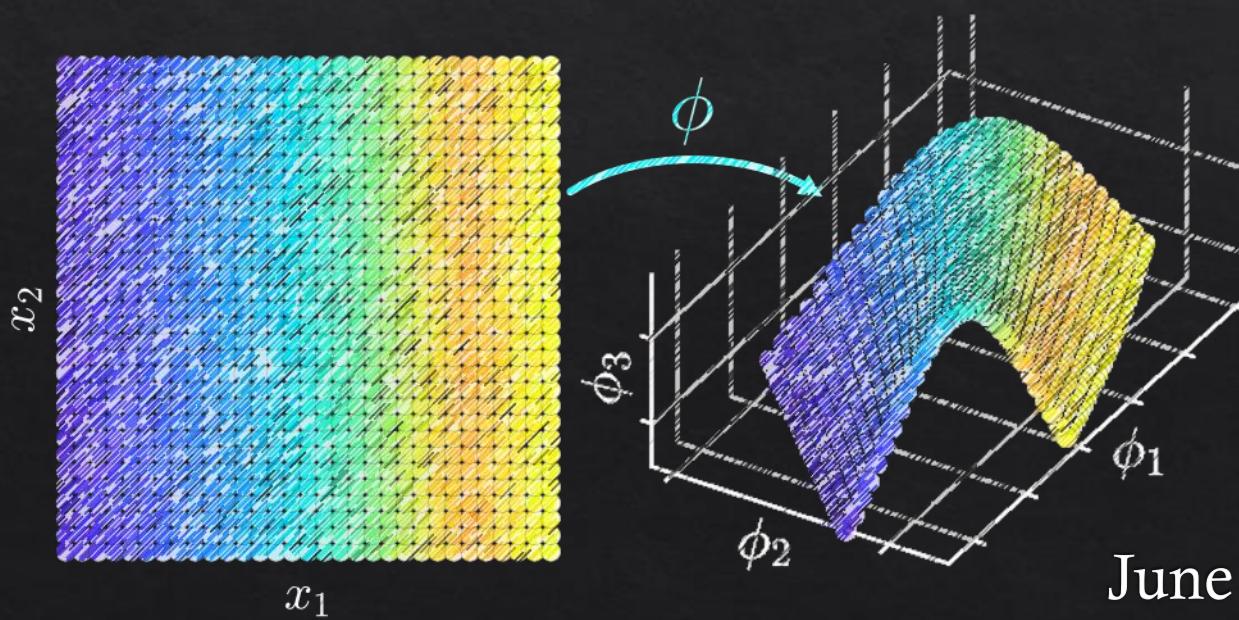


Supervised Learning – Nonlinear Classification

Machine Learning Methods – Lecture 5



June 2021

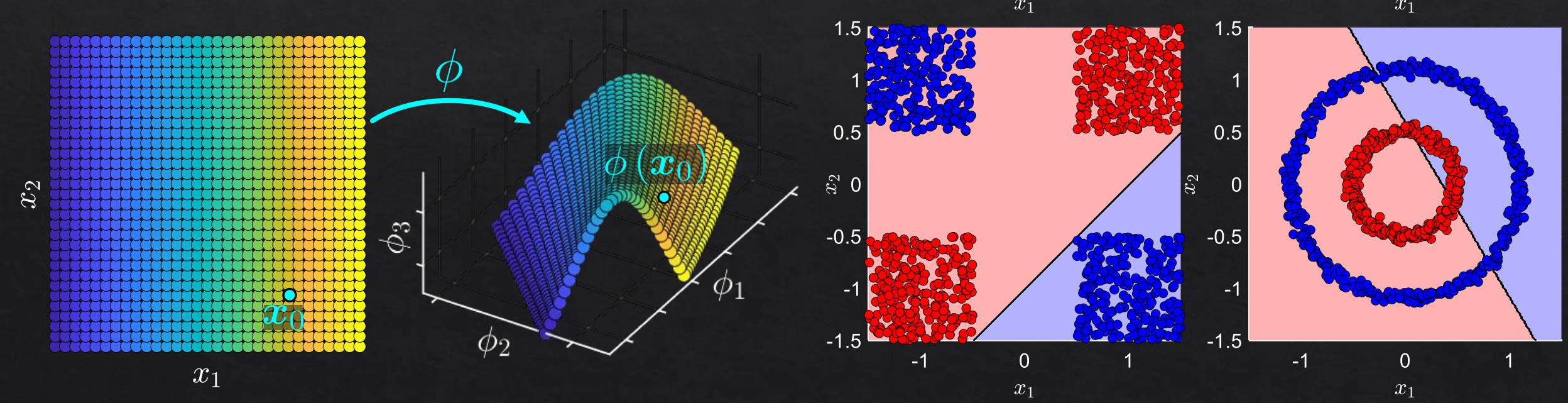
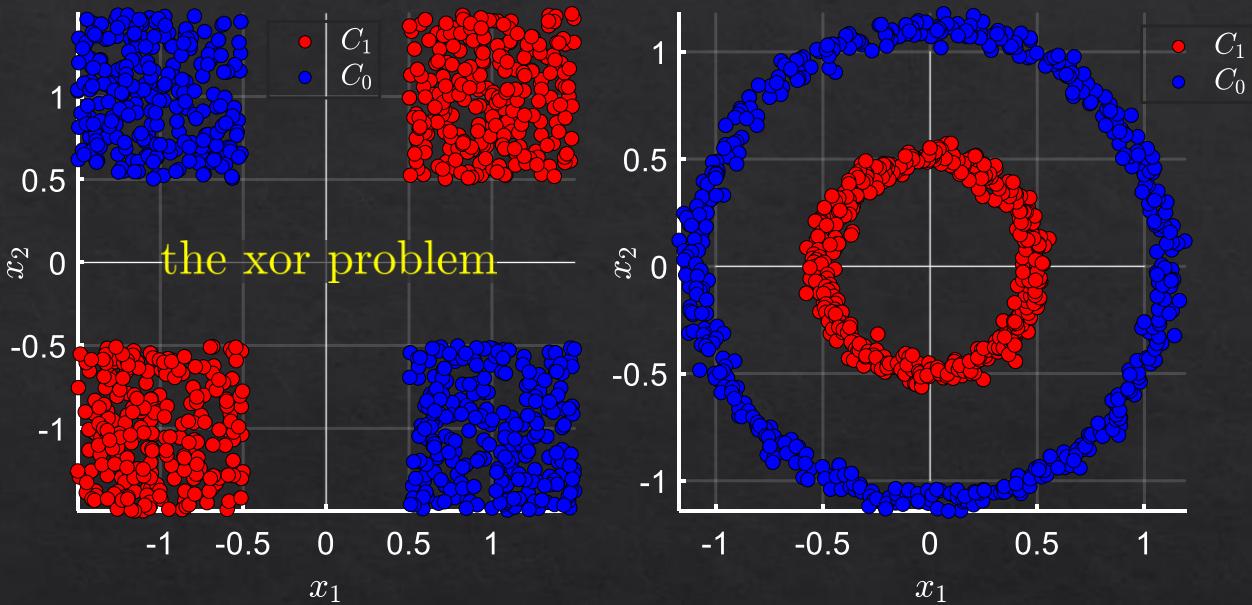


Linearly Non-separable Problems

- Consider the following classification problems:
- These problems are linearly non-separable in the space of (x_1, x_2) .

Can we find a mapping $(x_1, x_2) \mapsto \phi(x_1, x_2)$ such that non separable problem will turn into a linearly separable problem in the new space?

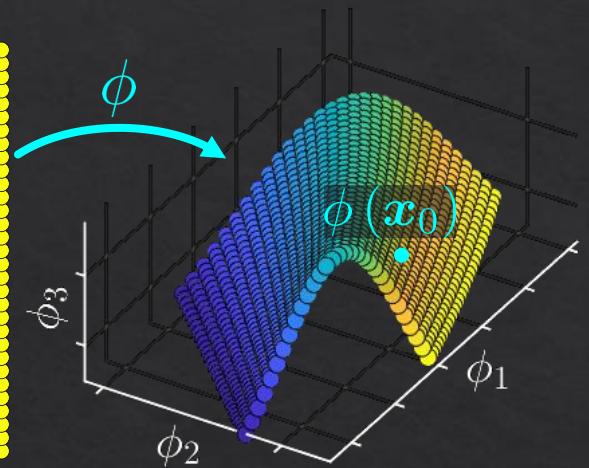
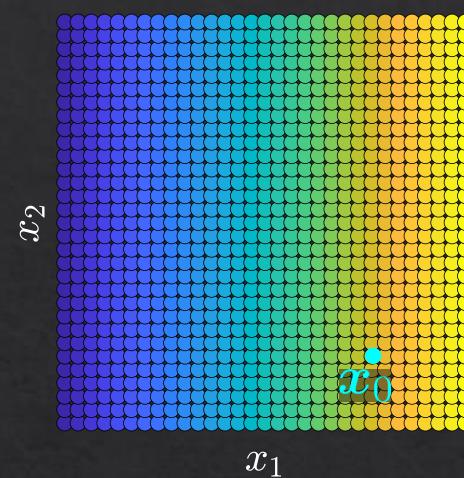
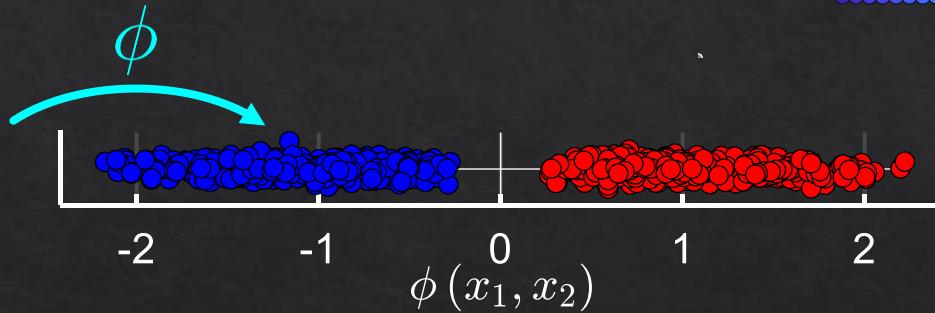
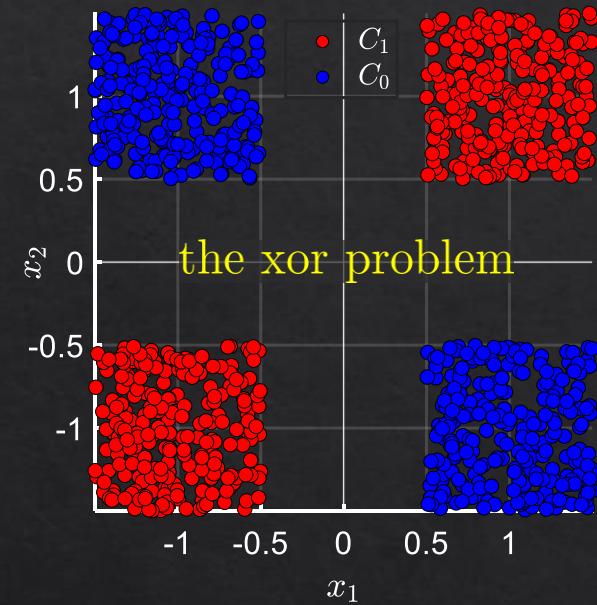
Yes! (and always!)



Feature Transform

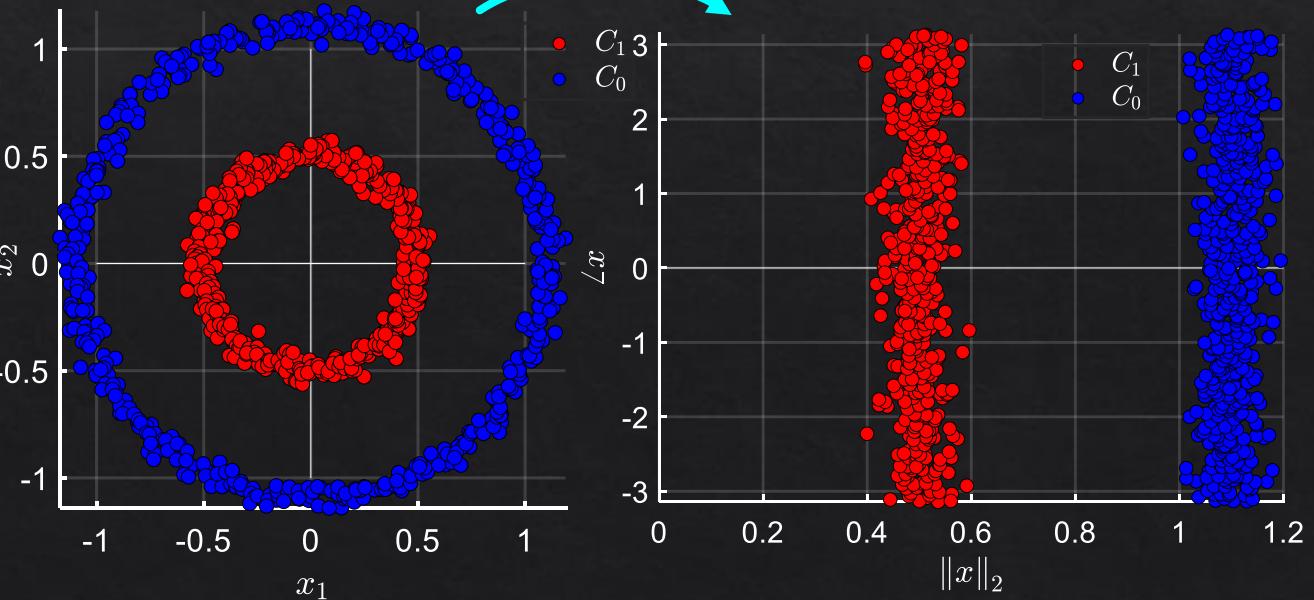
- Consider the following transformation:

$$\phi(x_1, x_2) = x_1 \cdot x_2$$



- Let:

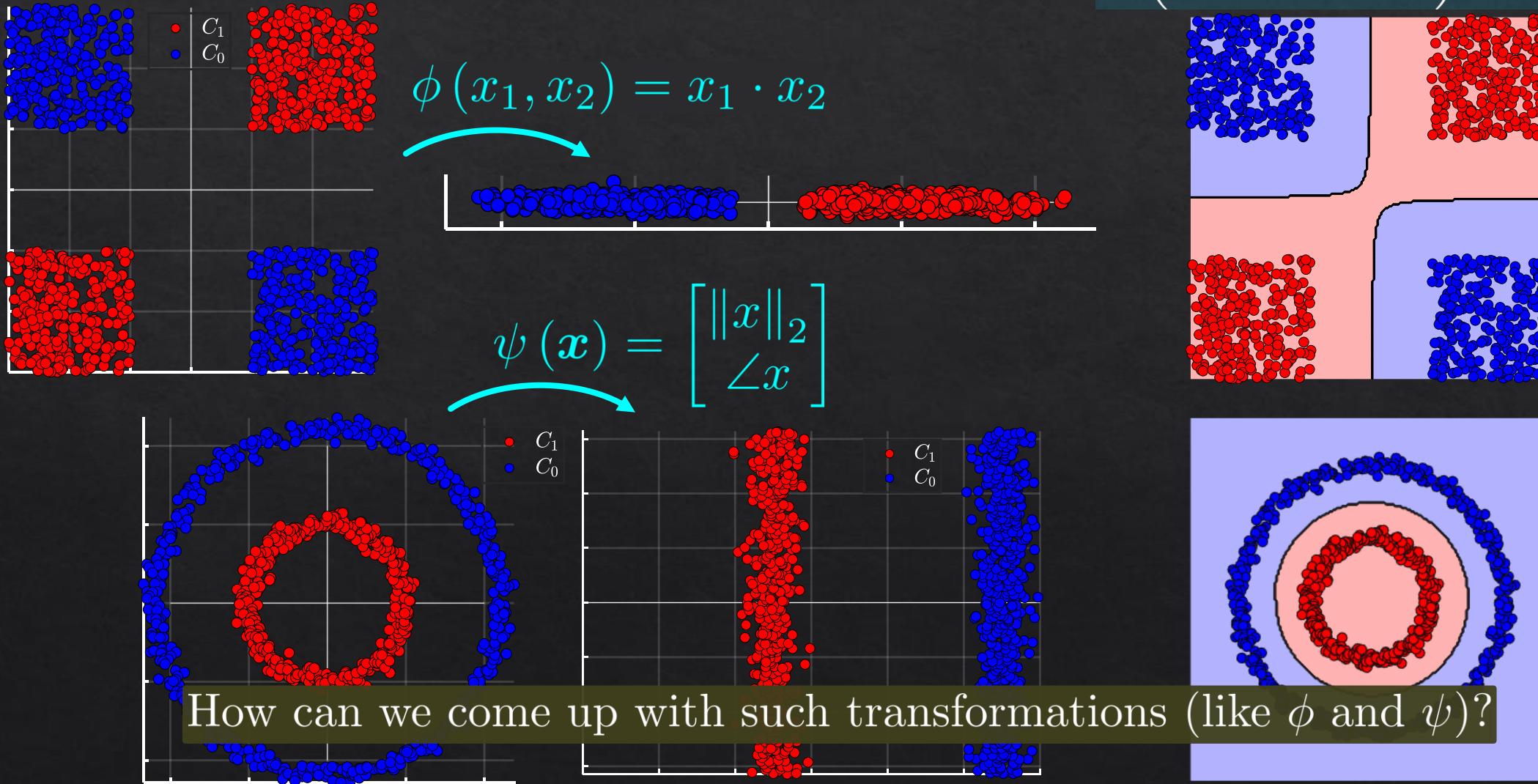
$$\psi(x_1, x_2) = \left(\underbrace{\sqrt{x_1^2 + x_2^2}}_{\|x\|_2}, \underbrace{\arctan\left(\frac{x_2}{x_1}\right)}_{\angle x} \right)$$



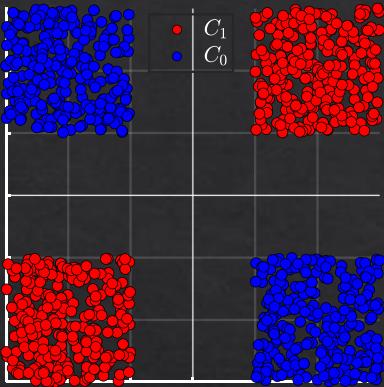
Feature Transform

- Instead of training a linear classifier in the original domain: $y_i (\mathbf{w}^T \mathbf{x}_i - b) > 0$

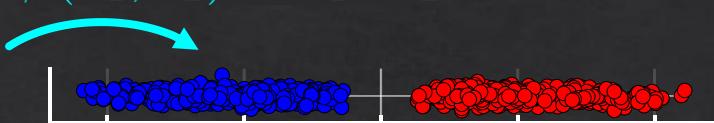
We can train a linear classifier in the new domain: $y_i (\tilde{\mathbf{w}}^T \phi(\mathbf{x}_i) - \tilde{b}) > 0$



The Kernel Trick – Outline

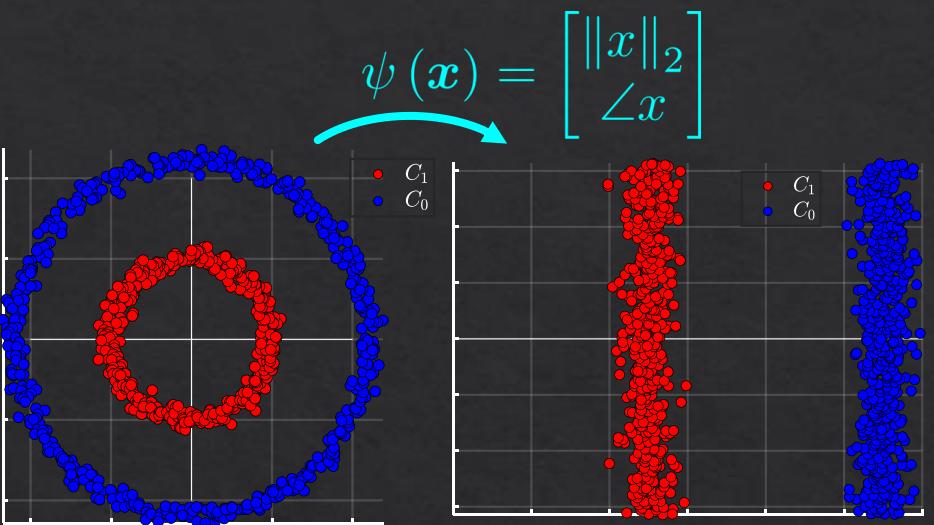


$$\phi(x_1, x_2) = x_1 \cdot x_2$$



Feature transform:

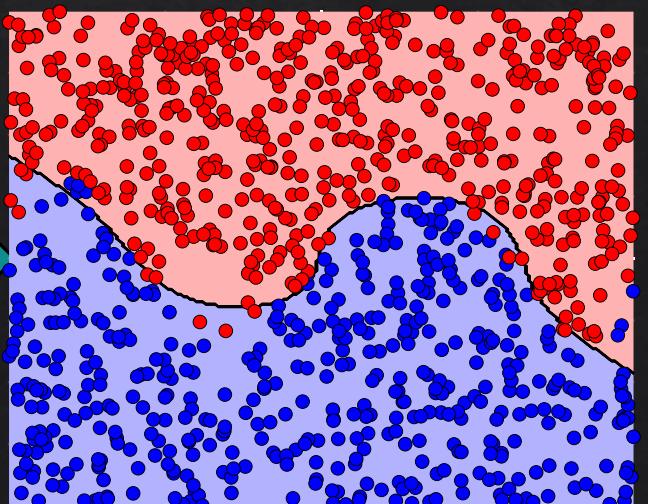
$$\{(\mathbf{x}_i, y_i)\}_{i=1}^N \longrightarrow \{(\phi(\mathbf{x}_i), y_i)\}_{i=1}^N$$



- To obtain a useful transformation, we now review:
 - Optimization dual problem
 - SVM dual problem
 - Kernel functions
 - The kernel trick



$$K(x_i, x_j) = (1 + x_i^T x_j)^3$$



The Dual Problem – Intuition

- Consider the following constraint optimization problem:

$$\begin{cases} \min_{x \in \mathbb{R}} x^2 \\ \text{subject to } x + 1 \leq 0 \end{cases} \xrightarrow{\text{same global minimum}} F(x) := \max_{\lambda \geq 0} x^2 + \lambda(x + 1) = \begin{cases} x^2 & (x + 1) \leq 0 \\ \infty & (x + 1) \geq 0 \end{cases}$$

- We can formulate this (primal) problem as follows:

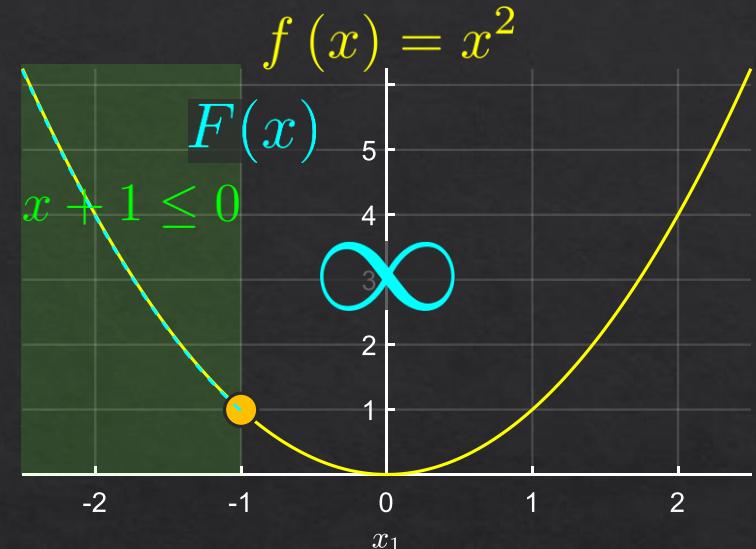
$$\min_{x \in \mathbb{R}} F(x) = \min_{x \in \mathbb{R}} \max_{\lambda \geq 0} \underbrace{x^2 + \lambda(x + 1)}_{=: L(x, \lambda)}$$

- The dual problem is defined by:

$$\max_{\lambda \geq 0} \min_{x \in \mathbb{R}} L(x, \lambda) = \max_{\lambda \geq 0} L(x^*, \lambda) = \max_{\lambda \geq 0} -\frac{\lambda^2}{4} + \lambda$$

primal dual $\lambda^* = 2$ $\implies x^* = -1$

$\underbrace{\min_{x \in \mathbb{R}} \max_{\lambda \geq 0} L(x, \lambda)}_{=} = \underbrace{\max_{\lambda \geq 0} \min_{x \in \mathbb{R}} L(x, \lambda)}_{=}$ strong duality (under some conditions)
 we can solve the dual instead of the primal



SVM – The Dual Problem

$$\min_{x \in \mathbb{R}} \max_{\lambda \geq 0} L(x, \lambda) = \max_{\lambda \geq 0} \min_{x \in \mathbb{R}} L(x, \lambda)$$

The primal SVM problem:

$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \\ \xi_i = \max \{0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i - b)\} \end{cases}$$

The dual SVM problem:

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{subject to} \\ 0 \leq \alpha_i \leq C \\ \sum_i \alpha_i y_i = 0 \end{cases}$$

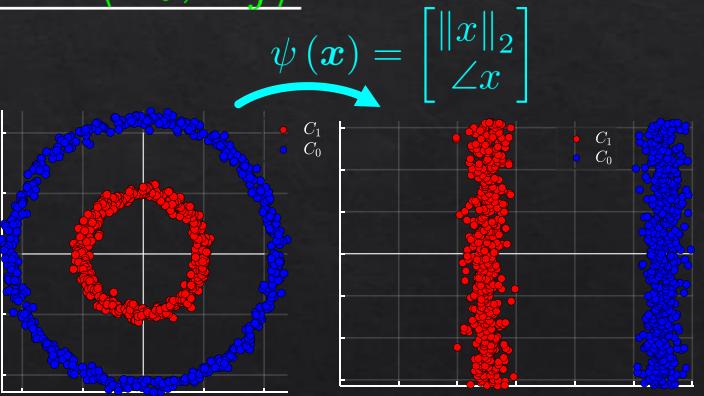
- Notice that the dual problem requires only the inner products $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$.

Feature transform:

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^N \rightarrow \{(\phi(\mathbf{x}_i), y_i)\}_{i=1}^N$$

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \\ \text{subject to} \\ 0 \leq \alpha_i \leq C \\ \sum_i \alpha_i y_i = 0 \end{cases}$$

using kernel functions we can compute the inner product $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$
without knowing $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$



Kernel Function – Example

- Let $x_i, x_j \in \mathbb{R}$:
- Consider the following kernel:

$$\begin{aligned} K(x_i, x_j) &= (1 + x_i x_j)^2 \\ &= 1 + 2x_i x_j + x_i^2 x_j^2 \end{aligned}$$

$$= \left\langle \begin{bmatrix} 1 \\ \sqrt{2}x_i \\ x_i^2 \end{bmatrix}, \begin{bmatrix} 1 \\ \sqrt{2}x_j \\ x_j^2 \end{bmatrix} \right\rangle$$
$$\implies K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \phi(x_i), \phi(x_j) \rangle \\ \text{subject to} \\ 0 \leq \alpha_i \leq C \\ \sum_i \alpha_i y_i = 0 \end{cases}$$

using kernel functions we can compute $\langle \phi(x_i), \phi(x_j) \rangle$ without knowing $\phi(x_i)$ and $\phi(x_j)$

$$\phi(x) := \begin{bmatrix} 1 \\ \sqrt{2}x \\ x^2 \end{bmatrix} \in \mathbb{R}^3$$

- $K(x_i, x_j)$ is a kernel because we can write it as an inner product:

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle \quad \text{for some (possibly unknown) } \phi.$$

Kernel Function

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underbrace{\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle}_{K(\mathbf{x}_i, \mathbf{x}_j)} \\ \text{subject to} \\ 0 \leq \alpha_i \leq C \\ \sum_i \alpha_i y_i = 0 \end{cases}$$

- A kernel function $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies:
 1. Symmetry: $K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_j, \mathbf{x}_i)$.
 2. K is positive (semi) definite.
- A kernel function K can be written as:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^M \phi_m(\mathbf{x}_i) \phi_m(\mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

$$\phi(\mathbf{x}_i) = \begin{bmatrix} \phi_1(\mathbf{x}_i) \\ \phi_2(\mathbf{x}_i) \\ \vdots \\ \phi_M(\mathbf{x}_i) \end{bmatrix} \quad M \text{ might be infinity}$$

Useful kernels

- Polynomial kernel ($d \geq 1$):

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^d$$

- Gaussian kernel ($\sigma > 0$):

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

- Laplacian kernel ($\sigma > 0$):

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|}{\sigma}\right)$$

- Sigmoid kernel:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh\left(\frac{1}{d} \mathbf{x}_i^T \mathbf{x}_j + c\right)$$

The Kernel Trick – Example

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$

- As a simple example, let $x_i, x_j \in \mathbb{R}$.
- Consider the following kernel:

$$\begin{aligned} K(x_i, x_j) &= (1 + x_i x_j)^2 \\ &= 1 + 2x_i x_j + x_i^2 x_j^2 \\ &= \left\langle \begin{bmatrix} 1 \\ \sqrt{2}x_i \\ x_i^2 \end{bmatrix}, \begin{bmatrix} 1 \\ \sqrt{2}x_j \\ x_j^2 \end{bmatrix} \right\rangle \\ \phi(x) &:= \begin{bmatrix} 1 \\ \sqrt{2}x \\ x^2 \end{bmatrix} \in \mathbb{R}^3 \end{aligned}$$

$$\implies K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$

- The Gaussian kernel:

$$K(x_i, x_j) = \exp\left(-\frac{(x_i - x_j)^2}{2\sigma^2}\right)$$

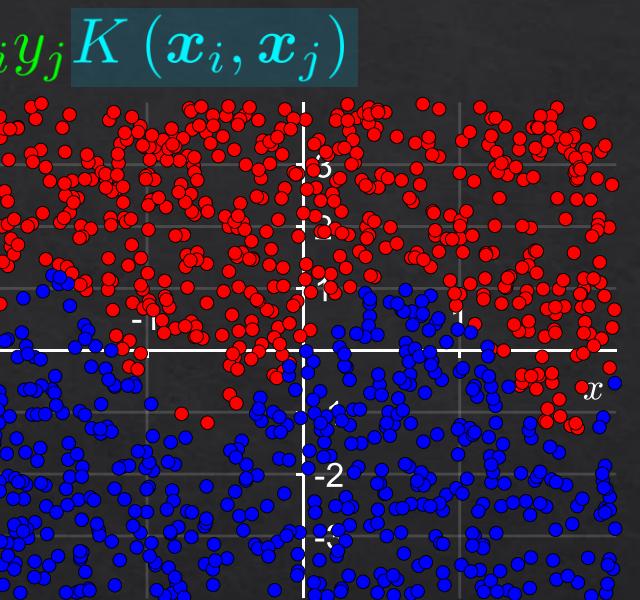
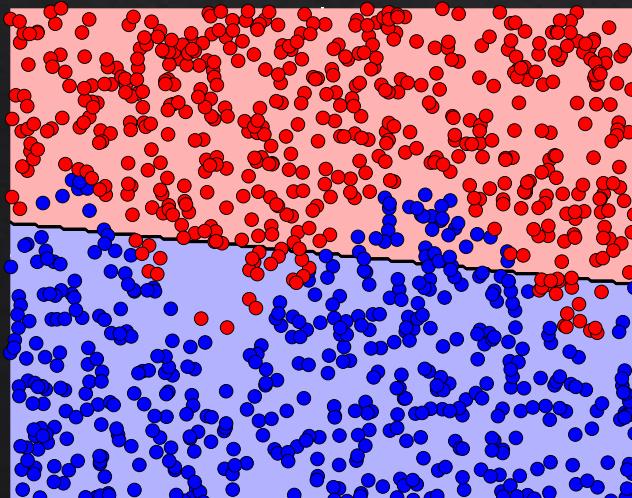
$$\phi(x) = \exp\left(-\frac{x^2}{2\sigma^2}\right) \cdot \begin{bmatrix} 1 \\ \frac{x}{\sigma} \\ \frac{1}{\sqrt{2!}} \frac{x^2}{\sigma^2} \\ \frac{1}{\sqrt{3!}} \frac{x^3}{\sigma^3} \\ \vdots \end{bmatrix} \in \mathbb{R}^\infty$$

The Kernel Trick – Polynomial Kernel

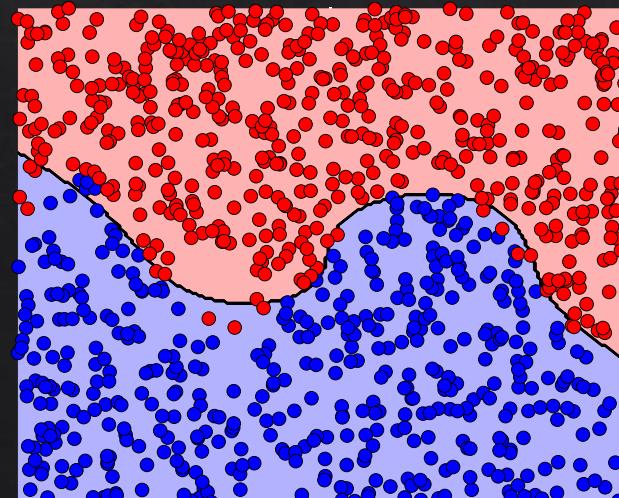
SVM dual problem:

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{subject to} \\ 0 \leq \alpha_i \leq C \\ \sum_i \alpha_i y_i = 0 \end{cases}$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^1$$



$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^3$$



$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

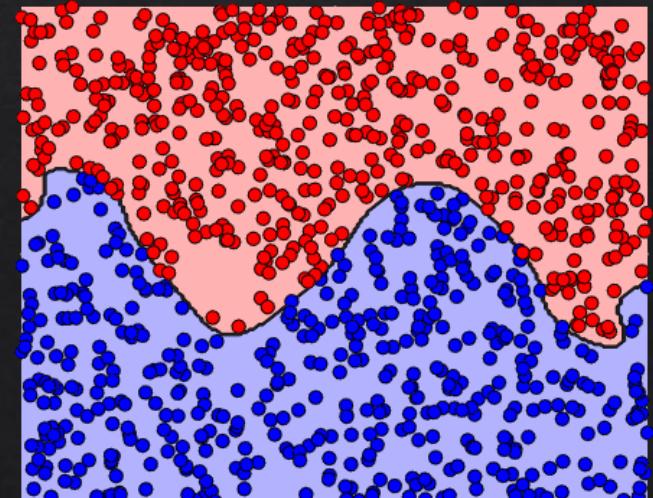
- Polynomial kernel ($d \geq 1$):

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^d$$

- Gaussian kernel ($\sigma > 0$):

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^5$$



Kernel SVM

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$

- Polynomial kernel ($d \geq 1$):

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^d$$

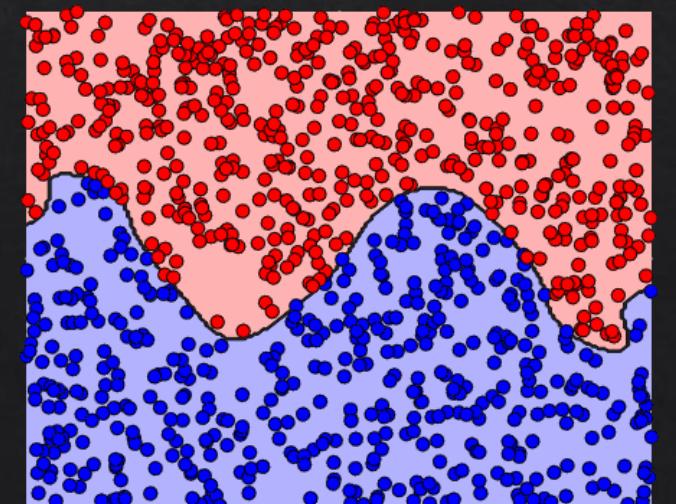
- Gaussian kernel ($\sigma > 0$):

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$



SVM with a Gaussian Kernel

$$K(x_i, x_j) = (1 + x_i^T x_j)^5$$

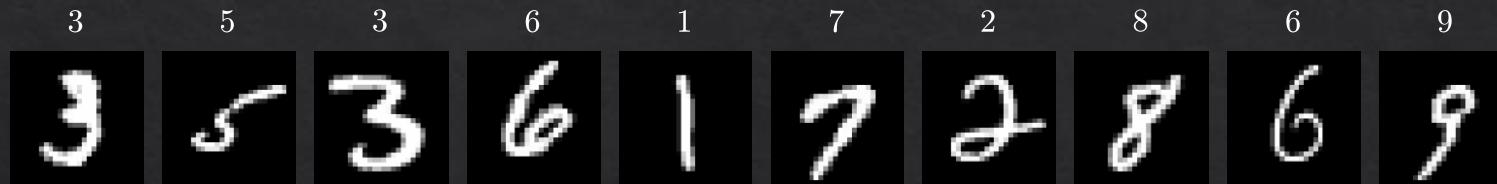


The Kernel Trick – Gaussian Kernel

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$



- Polynomial kernel ($d \geq 1$):
$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^d$$



- Gaussian kernel ($\sigma > 0$):
$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

Overfit

Fit

Underfit

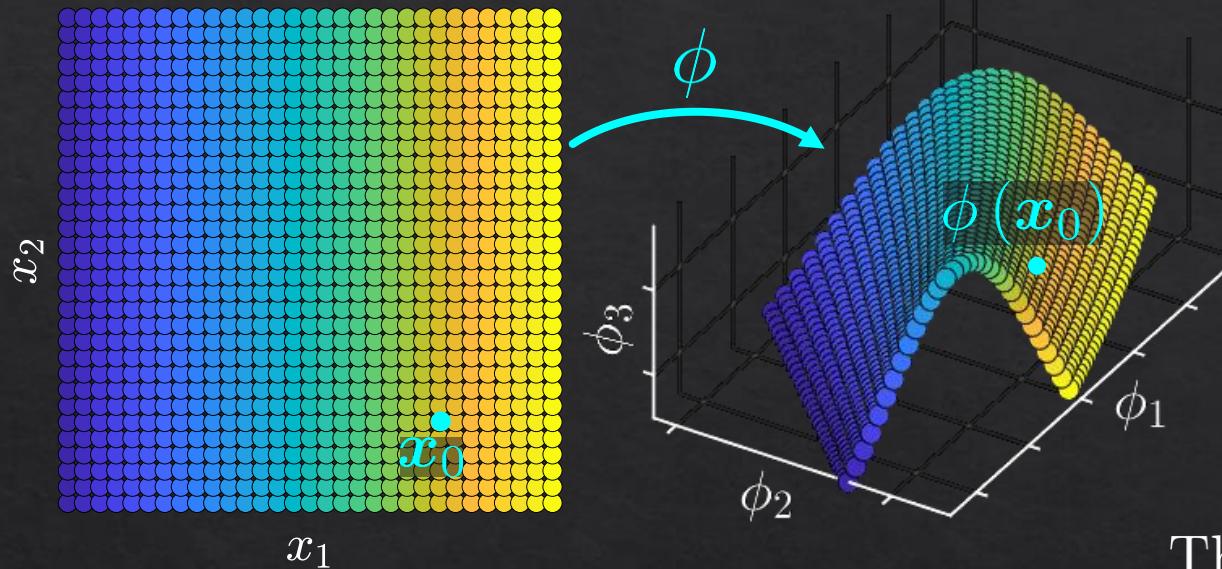
minimal
test error

σ

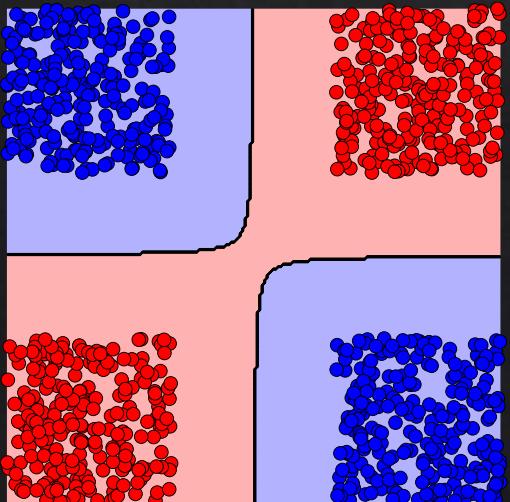


Questions

Feature transform:



Non-linear classification:



The kernel trick:

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^5$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

