

Kernel Methods

Machine Learning

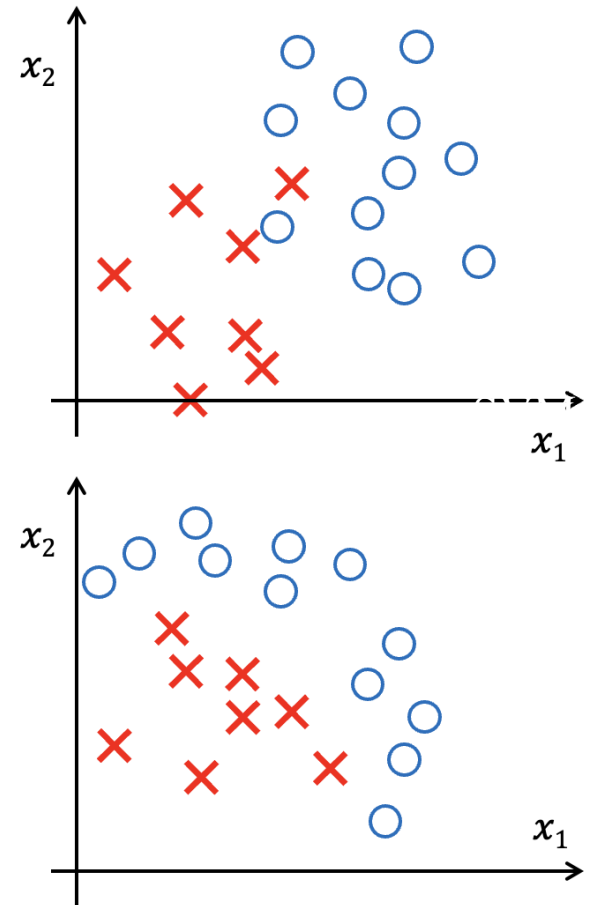
Hamid R Rabiee – Zahra Dehghanian
Spring 2025



Sharif University
of Technology

Not linearly separable data

- Noisy data or overlapping classes (we discussed about it: soft margin)
 - Near linearly separable
- Non-linear decision surface
 - Transform to a new feature space



Nonlinear SVM

- Assume a transformation $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$ on the feature space

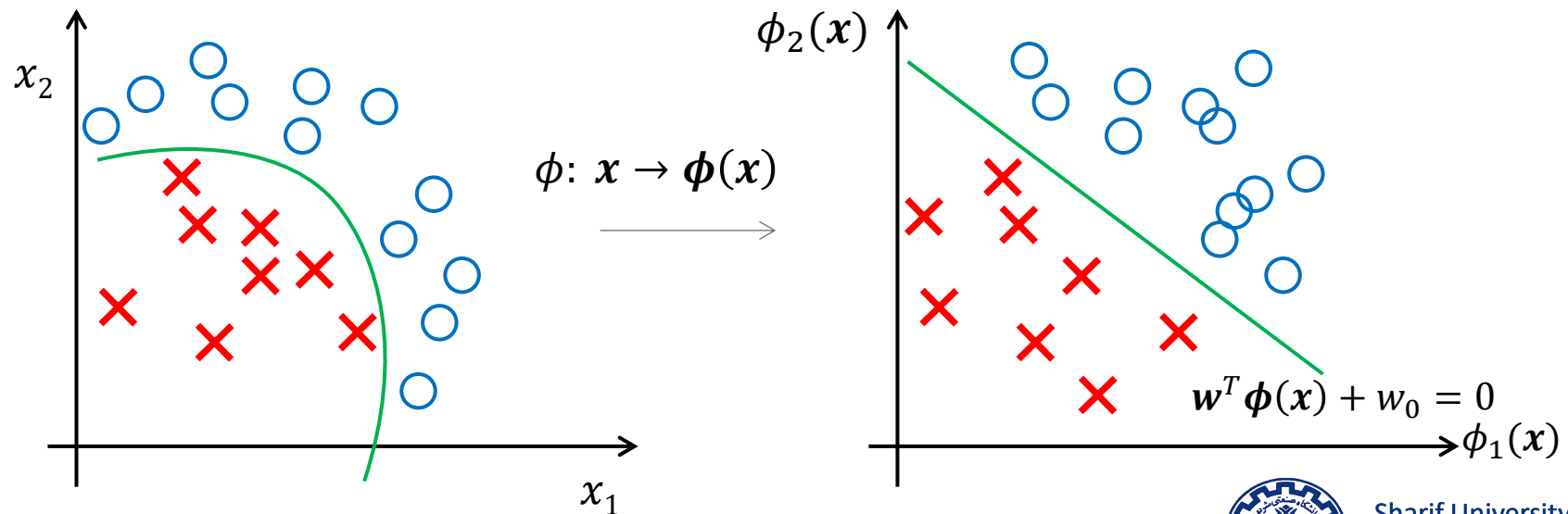
- $x \rightarrow \phi(x)$

$$\phi(x) = [\phi_1(x), \dots, \phi_m(x)]$$

$\{\phi_1(x), \dots, \phi_m(x)\}$: set of basis functions (or features)

$$\phi_i(x): \mathbb{R}^d \rightarrow \mathbb{R}$$

- Find a hyper-plane in the transformed feature space:



Soft-margin SVM in a transformed space: Primal problem

- Primal problem:

$$\begin{aligned} \min_{\mathbf{w}, w_0} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ \text{s. t.} \quad & y^{(n)} (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(n)}) + w_0) \geq 1 - \xi_n \quad n = 1, \dots, N \\ & \xi_n \geq 0 \end{aligned}$$

- $\mathbf{w} \in \mathbb{R}^m$: the weights that must be found
- If $m \gg d$ (very high dimensional feature space) then there are many more parameters to learn

Soft-margin SVM in a transformed space: Dual problem

- Optimization problem:

$$\max_{\alpha} \left\{ \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y^{(n)} y^{(m)} \boldsymbol{\phi}(\mathbf{x}^{(n)})^T \boldsymbol{\phi}(\mathbf{x}^{(m)}) \right\}$$

$$\text{Subject to } \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

$$0 \leq \alpha_n \leq C \quad n = 1, \dots, N$$

- If we have inner products $\boldsymbol{\phi}(\mathbf{x}^{(i)})^T \boldsymbol{\phi}(\mathbf{x}^{(j)})$, only $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]$ needs to be learnt.
 - not necessary to learn m parameters as opposed to the primal problem

Classifying a new data

$$\hat{y} = \text{sign}(w_0 + \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}))$$

$$\text{where } \mathbf{w} = \sum_{\alpha_n > 0} \alpha_n y^{(n)} \boldsymbol{\phi}(\mathbf{x}^{(n)})$$

$$\text{and } w_0 = y^{(s)} - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(s)})$$

Kernel SVM

- Learns linear decision boundary in a high dimension space without explicitly working on the mapped data
- Let $\phi(\mathbf{x})^T \phi(\mathbf{x}') = K(\mathbf{x}, \mathbf{x}')$ (kernel)
- Example: $\mathbf{x} = [x_1, x_2]$ and second-order ϕ :
$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$$K(\mathbf{x}, \mathbf{x}') = 1 + x_1 x_1' + x_2 x_2' + x_1^2 x_1'^2 + x_2^2 x_2'^2 + x_1 x_1' x_2 x_2'$$

Kernel trick

- Compute $K(x, x')$ without transforming x and x'
- Example: Consider $K(x, x') = (1 + x^T x')^2$
$$= (1 + x_1 x'_1 + x_2 x'_2)^2$$
$$= 1 + 2x_1 x'_1 + 2x_2 x'_2 + x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1 x'_1 x_2 x'_2$$

This is an inner product in:

$$\phi(x) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2]$$
$$\phi(x') = [1, \sqrt{2}x'_1, \sqrt{2}x'_2, x_1'^2, x_2'^2, \sqrt{2}x'_1 x'_2]$$

Polynomial kernel: Degree two

- We instead use $K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^2$ that corresponds to:

d -dimensional feature space $\mathbf{x} = [x_1, \dots, x_d]^T$

$\phi(\mathbf{x})$

$$= [1, \sqrt{2}x_1, \dots, \sqrt{2}x_d, x_1^2, \dots, x_d^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_1x_d, \sqrt{2}x_2x_3, \dots, \sqrt{2}x_{d-1}x_d]^T$$

Polynomial kernel

- This can similarly be generalized to d-dimension \mathbf{x} and ϕ s are polynomials of order M :

$$\begin{aligned}K(\mathbf{x}, \mathbf{x}') &= (1 + \mathbf{x}^T \mathbf{x}')^M \\&= (1 + x_1 x'_1 + x_2 x'_2 + \cdots + x_d x'_d)^M\end{aligned}$$

- Example: SVM boundary for a polynomial kernel

- $w_0 + \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) = 0$

$$\Rightarrow w_0 + \sum_{\alpha_i > 0} \alpha_i y^{(i)} \boldsymbol{\phi}(\mathbf{x}^{(i)})^T \boldsymbol{\phi}(\mathbf{x}) = 0$$

$$\Rightarrow w_0 + \sum_{\alpha_i > 0} \alpha_i y^{(i)} k(\mathbf{x}^{(i)}, \mathbf{x}) = 0$$

$$\Rightarrow w_0 + \sum_{\alpha_i > 0} \alpha_i y^{(i)} (1 + \mathbf{x}^{(i)T} \mathbf{x})^M = 0$$

➡ Boundary is a polynomial of order M

Why kernel?

- kernel functions K can indeed be efficiently computed, with a cost proportional to d (the dimensionality of the input) instead of m .
- Example: consider the second-order polynomial transform:

$$\boldsymbol{\phi}(\mathbf{x}) = [1, x_1, \dots, x_d, x_1^2, x_1x_2, \dots, x_dx_d]^T \quad q = 1 + d + d^2$$

$$\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}') = 1 + \sum_{i=1}^d x_i x'_i + \underbrace{\sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j}_{\sum_{i=1}^d x_i x'_i \times \sum_{j=1}^d x_j x'_j} \quad O(q)$$

$$\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}') = 1 + (\mathbf{x}^T \mathbf{x}') + (\mathbf{x}^T \mathbf{x}')^2$$

$O(d)$

Gaussian or RBF kernel

- If $K(\mathbf{x}, \mathbf{x}')$ is an inner product in some transformed space of \mathbf{x} , it is good

- $K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma}\right)$

- Take one dimensional case with $\gamma = 1$:

$$K(x, x') = \exp(-(x - x')^2)$$

$$= \exp(-x^2) \exp(-x'^2) \sum_{k=0}^{\infty} \frac{2^k x^k x'^k}{k!}$$

Some common kernel functions

- Linear: $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$
- Polynomial: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^M$
- Gaussian: $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma}\right)$
- Sigmoid: $k(\mathbf{x}, \mathbf{x}') = \tanh(a\mathbf{x}^T \mathbf{x}' + b)$

Kernel formulation of SVM

- Optimization problem:

$$\max_{\alpha} \left\{ \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y^{(n)} y^{(m)} k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \right\}$$

$$\text{Subject to } \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

$$0 \leq \alpha_n \leq C \quad n = 1, \dots, N$$

$$Q = \begin{bmatrix} y^{(1)}y^{(1)}K(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & y^{(1)}y^{(N)}K(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) \\ \vdots & \ddots & \vdots \\ y^{(N)}y^{(1)}K(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) & \dots & y^{(N)}y^{(N)}K(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix}$$

Classifying a new data

$$\hat{y} = \text{sign}(w_0 + \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}))$$

$$\text{where } \mathbf{w} = \sum_{\alpha_n > 0} \alpha_n y^{(n)} \boldsymbol{\phi}(\mathbf{x}^{(n)})$$

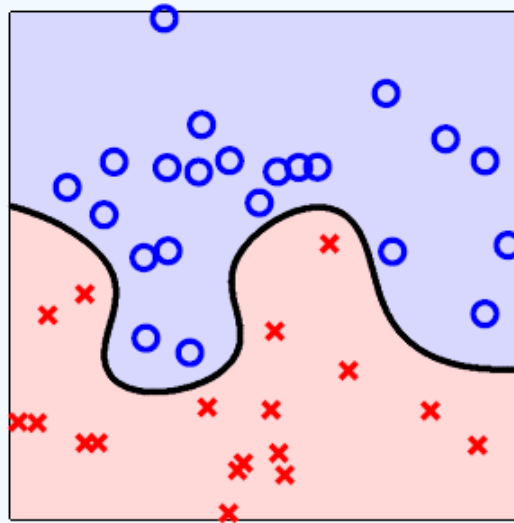
$$\text{and } w_0 = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(s)})$$

$$\hat{y} = \text{sign} \left(w_0 + \sum_{\alpha_n > 0} \alpha_n y^{(n)} k(\mathbf{x}^{(n)}, \mathbf{x}) \right)$$
$$w_0 = y^{(s)} - \sum_{\alpha_n > 0} \alpha_n y^{(n)} k(\mathbf{x}^{(n)}, \mathbf{x}^{(s)})$$

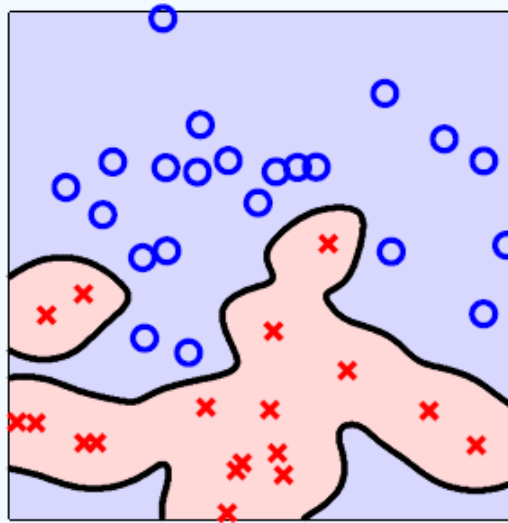
Gaussian kernel

- Example: SVM boundary for a Gaussian kernel
 - Considers a Gaussian function around each data point.
 - $w_0 + \sum_{\alpha_i > 0} \alpha_i y^{(i)} \exp\left(-\frac{\|x - x^{(i)}\|^2}{\sigma}\right) = 0$
 - SVM + Gaussian Kernel can classify any arbitrary training set
 - Training error is zero when $\sigma \rightarrow 0$
 - All samples become support vectors (likely overfitting)

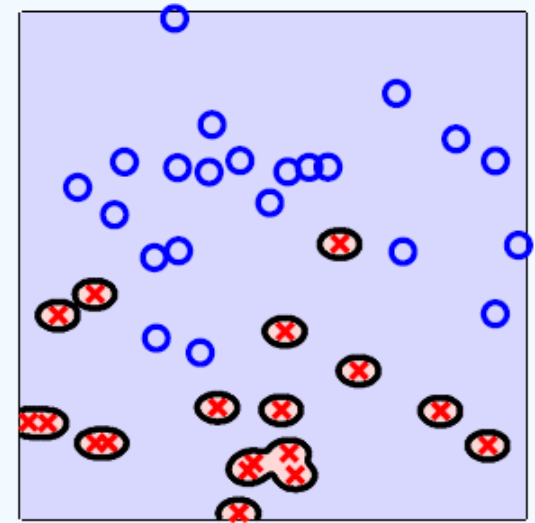
Hard margin Example



$$\exp(-1\|\mathbf{x} - \mathbf{x}'\|^2)$$



$$\exp(-10\|\mathbf{x} - \mathbf{x}'\|^2)$$



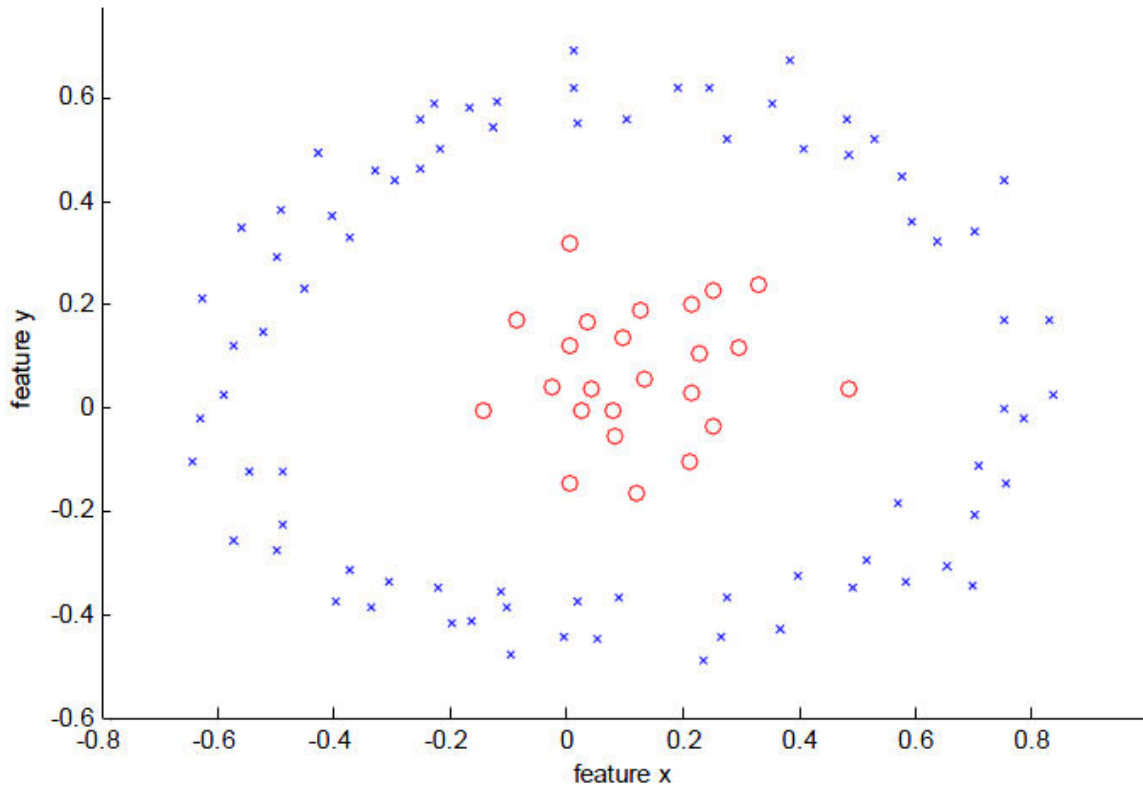
$$\exp(-100\|\mathbf{x} - \mathbf{x}'\|^2)$$

Y. Abu-Mostafa et. Al, 2012

- For narrow Gaussian (large σ), even the protection of a large margin cannot suppress overfitting.

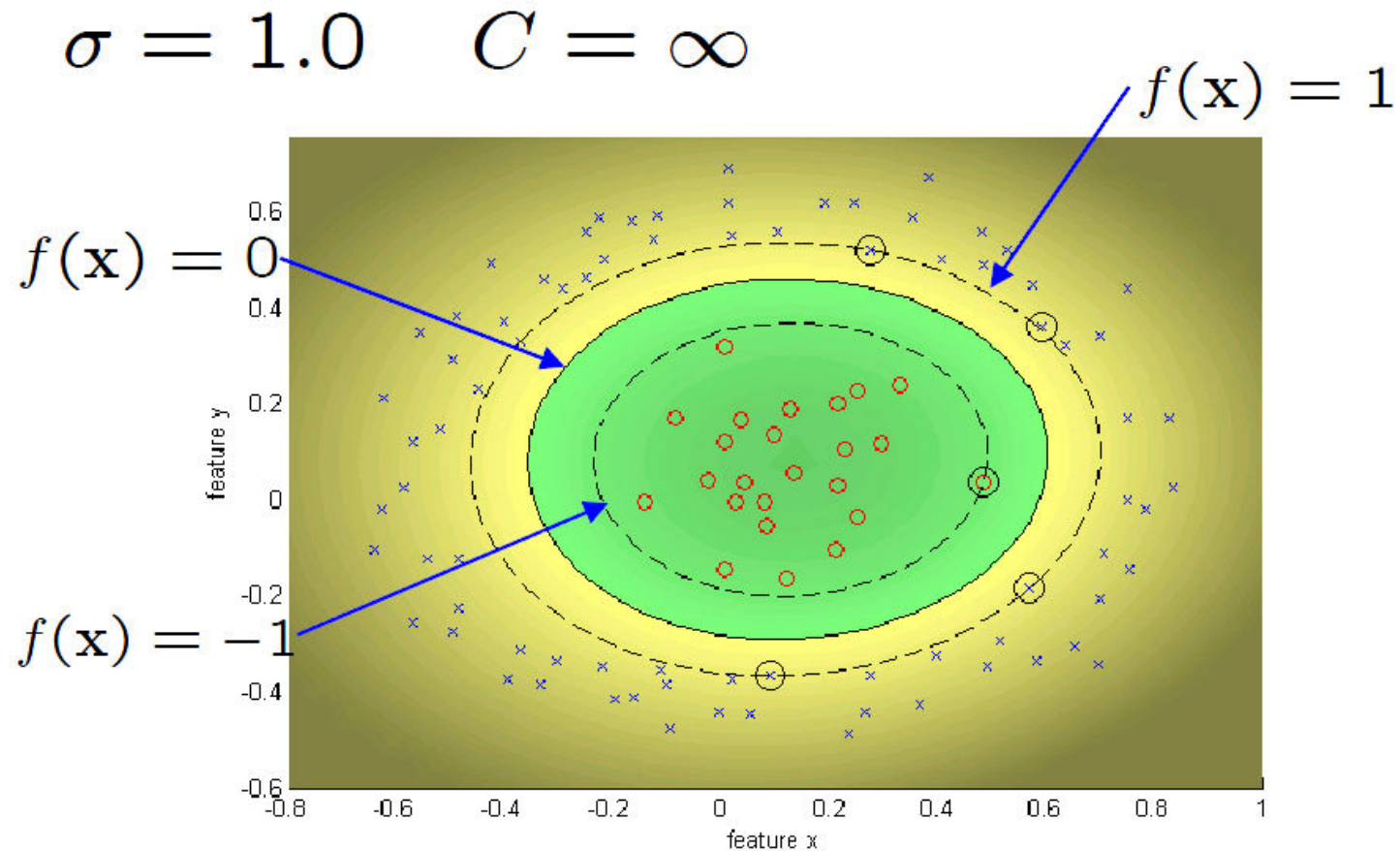
SVM Gaussian kernel: Example

$$f(\mathbf{x}) = w_0 + \sum_{\alpha_i > 0} \alpha_i y^{(i)} \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}^{(i)}\|^2}{2\sigma^2}\right)$$



This example has been adopted from Zisserman's slides

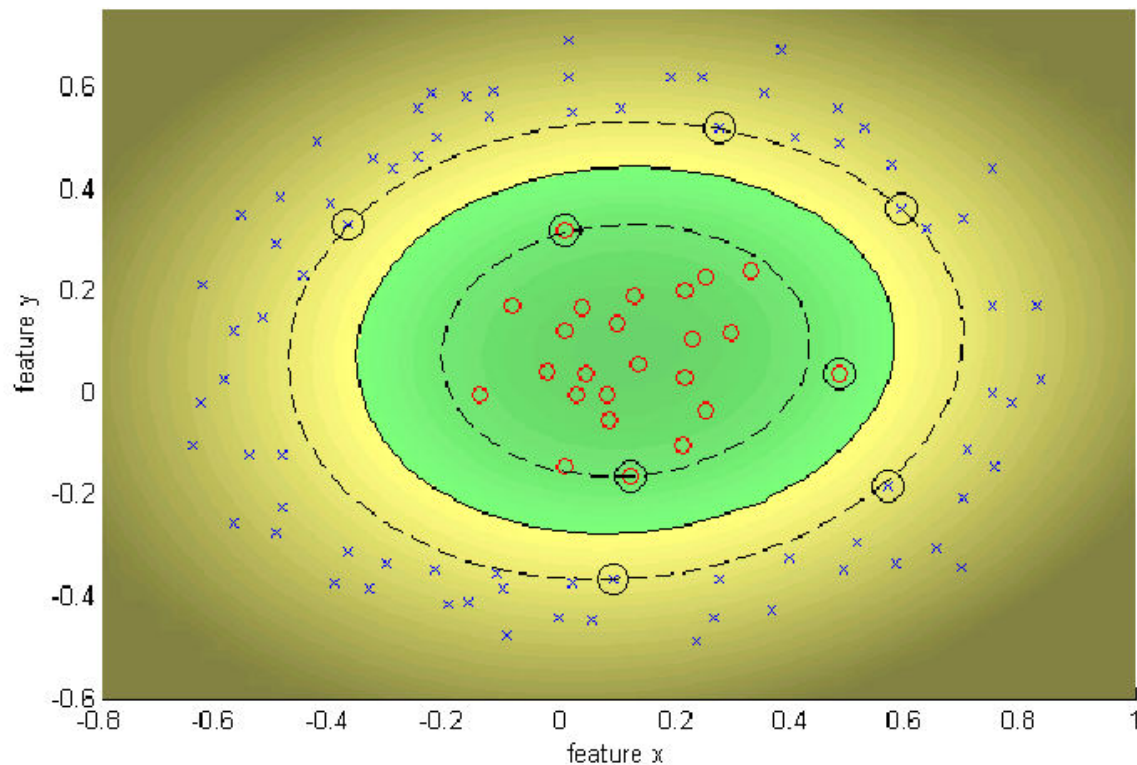
SVM Gaussian kernel: Example



This example has been adopted from Zisserman's slides

SVM Gaussian kernel: Example

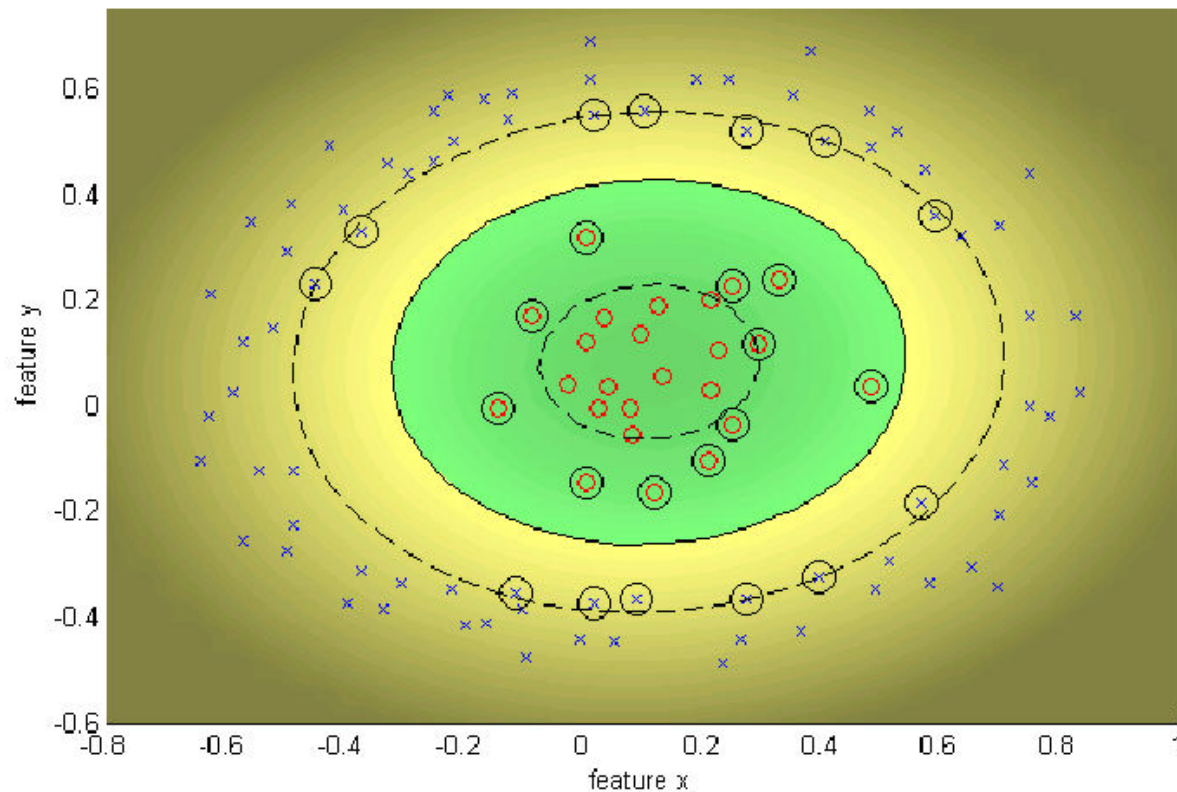
$$\sigma = 1.0 \quad C = 100$$



This example has been adopted from Zisserman's slides

SVM Gaussian kernel: Example

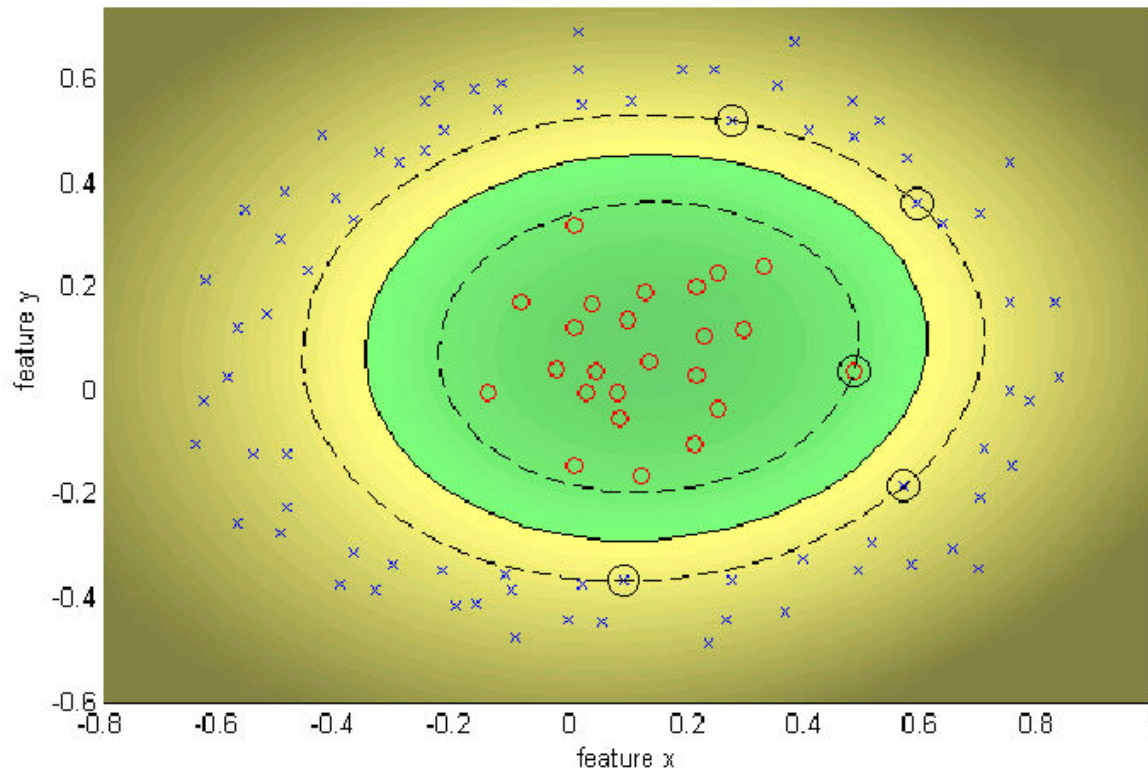
$$\sigma = 1.0 \quad C = 10$$



This example has been adopted from Zisserman's slides

SVM Gaussian kernel: Example

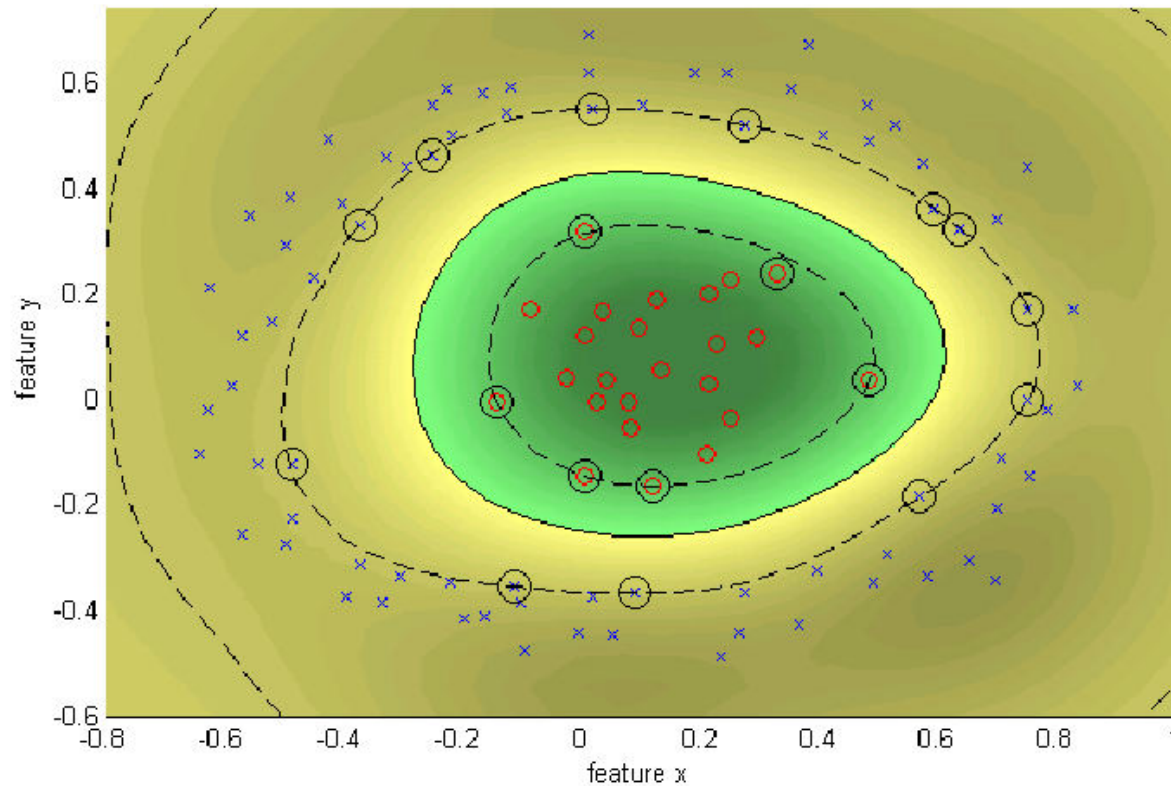
$$\sigma = 1.0 \quad C = \infty$$



This example has been adopted from Zisserman's slides

SVM Gaussian kernel: Example

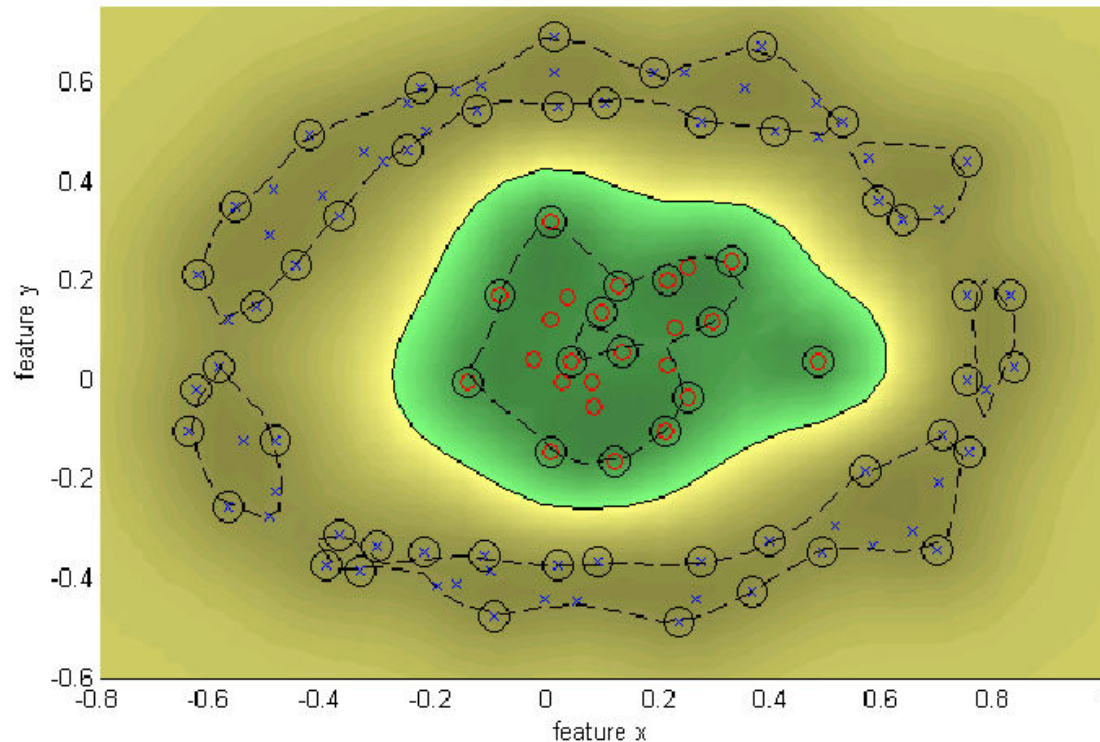
$$\sigma = 0.25 \quad C = \infty$$



This example has been adopted from Zisserman's slides

SVM Gaussian kernel: Example

$$\sigma = 0.1 \quad C = \infty$$



This example has been adopted from Zisserman's slides

Kernel trick: Idea

- Kernel trick → Extension of many well-known algorithms to kernel-based ones
 - By substituting the dot product with the kernel function
 - $k(x, x') = \phi(x)^T \phi(x')$
 - $k(x, x')$ shows the dot product of x and x' in the transformed space.
- Idea: when the input vectors appears only in the form of dot products, we can use kernel trick
 - Solving the problem without explicitly mapping the data
 - Explicit mapping is expensive if $\phi(x)$ is very high dimensional

Kernel trick: Idea (Cont'd)

- Instead of using a mapping $\phi: \mathcal{X} \leftarrow \mathcal{F}$ to represent $x \in \mathcal{X}$ by $\phi(x) \in \mathcal{F}$, a kernel function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is used.
 - We specify only an inner product function between points in the transformed space (not their coordinates)
 - In many cases, the inner product in the embedding space can be computed efficiently.

Constructing kernels

- Construct kernel functions directly
 - Ensure that it is a valid kernel
 - Corresponds to an inner product in some feature space.
- Example: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$
 - Corresponding mapping: $\boldsymbol{\phi}(\mathbf{x}) = [x_1^2, \sqrt{2}x_1x_2, x_2^2]^T$ for $\mathbf{x} = [x_1, x_2]^T$
- We need a way to test whether a kernel is valid without having to construct $\boldsymbol{\phi}(\mathbf{x})$

Construct Valid Kernels

$$\begin{aligned}k(\mathbf{x}, \mathbf{x}') &= ck_1(\mathbf{x}, \mathbf{x}') \\k(\mathbf{x}, \mathbf{x}') &= f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \\k(\mathbf{x}, \mathbf{x}') &= q(k_1(\mathbf{x}, \mathbf{x}')) \\k(\mathbf{x}, \mathbf{x}') &= \exp(k_1(\mathbf{x}, \mathbf{x}')) \\k(\mathbf{x}, \mathbf{x}') &= k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \\k(\mathbf{x}, \mathbf{x}') &= k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \\k(\mathbf{x}, \mathbf{x}') &= k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \\k(\mathbf{x}, \mathbf{x}') &= \mathbf{x}^T \mathbf{A} \mathbf{x}' \\k(\mathbf{x}, \mathbf{x}') &= k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \\k(\mathbf{x}, \mathbf{x}') &= k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b)\end{aligned}$$

- $c > 0$, k_1 : valid kernel
- $f(\cdot)$: any function
- $q(\cdot)$: a polynomial with coefficients ≥ 0
- k_1, k_2 : valid kernels
- $\phi(\mathbf{x})$: a function from \mathbf{x} to \mathbb{R}^M
 $k_3(\cdot, \cdot)$: a valid kernel in \mathbb{R}^M
- \mathbf{A} : a symmetric positive semi-definite matrix
- \mathbf{x}_a and \mathbf{x}_b are variables (not necessarily disjoint) with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, and k_a and k_b are valid kernel functions over their respective spaces.

[Bishop]

Valid kernel: Necessary & sufficient conditions

- Gram matrix $\mathbf{K}_{N \times N}$: $K_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$
 - Restricting the kernel function to a set of points $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix}$$

- **Mercer** Theorem: The kernel matrix is **Symmetric Positive Semi-Definite** (for any choice of data points)
 - Any symmetric positive definite matrix can be regarded as a kernel matrix, that is as an inner product matrix in some space

[Shawe-Taylor & Cristianini 2004]

Extending linear methods to kernelized ones

- Kernelized version of linear methods
 - Linear methods are famous
 - Unique optimal solutions, faster learning algorithms, and better analysis
 - However, we often require nonlinear methods in real-world problems and so we can use kernel-based version of these linear algorithms
- Replacing inner products with kernels in linear algorithms \Rightarrow very flexible methods
 - We can operate in the mapped space without ever computing the coordinates of the data in that space

Example: kernelized minimum distance classifier

- If $\|\mathbf{x} - \boldsymbol{\mu}_1\| < \|\mathbf{x} - \boldsymbol{\mu}_2\|$ then assign \mathbf{x} to \mathcal{C}_1

$$(\mathbf{x} - \boldsymbol{\mu}_1)^T (\mathbf{x} - \boldsymbol{\mu}_1) < (\mathbf{x} - \boldsymbol{\mu}_2)^T (\mathbf{x} - \boldsymbol{\mu}_2)$$

$$-2\mathbf{x}^T \boldsymbol{\mu}_1 + \boldsymbol{\mu}_1^T \boldsymbol{\mu}_1 < -2\mathbf{x}^T \boldsymbol{\mu}_2 + \boldsymbol{\mu}_2^T \boldsymbol{\mu}_2$$

$$-2 \frac{\sum_{y^{(n)}=1} \mathbf{x}^T \mathbf{x}^{(n)}}{N_1} + \frac{\sum_{y^{(n)}=1} \sum_{y^{(m)}=1} \mathbf{x}^{(n)T} \mathbf{x}^{(m)}}{N_1 \times N_1} < -2 \frac{\sum_{y^{(n)}=2} \mathbf{x}^T \mathbf{x}^{(n)}}{N_2} + \frac{\sum_{y^{(n)}=2} \sum_{y^{(m)}=2} \mathbf{x}^{(n)T} \mathbf{x}^{(m)}}{N_2 \times N_2}$$

$$-2 \frac{\sum_{y^{(n)}=1} K(\mathbf{x}, \mathbf{x}^{(n)})}{N_1} + \frac{\sum_{y^{(n)}=1} \sum_{y^{(m)}=1} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)})}{N_1 \times N_1} < -2 \frac{\sum_{y^{(n)}=2} K(\mathbf{x}, \mathbf{x}^{(n)})}{N_2} + \frac{\sum_{y^{(n)}=2} \sum_{y^{(m)}=2} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)})}{N_2 \times N_2}$$

Which information can be obtained from kernel?

- Example: we know all pairwise distances
 - $d(\phi(x), \phi(z))^2 = \|\phi(x) - \phi(z)\|^2 = k(x, x) + k(z, z) - 2k(x, z)$
 - Therefore, we also know distance of points from center of mass of a set
- Many dimensionality reduction, clustering, and classification methods can be described according to pairwise distances.
 - This allow us to introduce kernelized versions of them

Kernels for structured data

- Kernels also can be defined on general types of data
 - Kernel functions do not need to be defined over vectors
 - just we need a symmetric positive definite matrix
- Thus, many algorithms can work with general (non-vectorial) data
 - Kernels exist to embed strings, trees, graphs, ...
- This may be more important than nonlinearity
 - kernel-based version of classical learning algorithms for recognition of structured data

Kernel function for objects

- Sets: Example of kernel function for sets:

$$k(A, B) = 2^{|A \cap B|}$$

- Strings: The inner product of the feature vectors for two strings can be defined as
 - e.g., sum over all common subsequences weighted according to their frequency of occurrence and lengths

A	E	G	A	T	E	A	G	G				
E	G	T	E	A	G	A	E	G	A	T	G	

Kernel trick advantages: summary

- Operating in the **mapped space without ever computing** the coordinates of the data in that space
- Besides vectors, we can introduce kernel functions for **structured data** (graphs, strings, etc.)
- Much of the **geometry** of the data in the embedding space is contained in all pairwise dot products
- In many cases, inner product in the embedding space can be computed efficiently.

Resources

- C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 6.1-6.2, 7.1.
- Yaser S. Abu-Mostafa, et al., “Learning from Data”, Chapter 8.